

2

LABORATORY FOR
COMPUTER SCIENCE



MASSACHUSETTS
INSTITUTE OF
TECHNOLOGY

MIT/LCS/TM-442

DTIC FILE COPY

AD-A232 829

ARE WAIT-FREE ALGORITHMS FAST?

DTIC
ELECTE
MAR 27 1991
S D D

Hagit Attiya
Nancy Lynch
Nir Shavit

DISTRIBUTION STATEMENT A

Approved for public release
Distribution Unlimited

March 1991

545 TECHNOLOGY SQUARE, CAMBRIDGE, MASSACHUSETTS 02139

91 3 22 089

REPORT DOCUMENTATION PAGE

1a. REPORT SECURITY CLASSIFICATION Unclassified			1b. RESTRICTIVE MARKINGS		
2a. SECURITY CLASSIFICATION AUTHORITY			3. DISTRIBUTION/AVAILABILITY OF REPORT Approved for public release; distribution is unlimited.		
2b. DECLASSIFICATION/DOWNGRADING SCHEDULE					
4. PERFORMING ORGANIZATION REPORT NUMBER(S) MIT/LCS/TM 442			5. MONITORING ORGANIZATION REPORT NUMBER(S) N00014-89-J-1988		
6a. NAME OF PERFORMING ORGANIZATION MIT Lab for Computer Science		6b. OFFICE SYMBOL (If applicable)		7a. NAME OF MONITORING ORGANIZATION Office of Naval Research/Dept. of Navy	
6c. ADDRESS (City, State, and ZIP Code) 545 Technology Square Cambridge, MA 02139			7b. ADDRESS (City, State, and ZIP Code) Information Systems Program Arlington, VA 22217		
8a. NAME OF FUNDING/SPONSORING ORGANIZATION DARPA/DOD		8b. OFFICE SYMBOL (If applicable)		9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER	
8c. ADDRESS (City, State, and ZIP Code) 1400 Wilson Blvd. Arlington, VA 22217			10. SOURCE OF FUNDING NUMBERS		
			PROGRAM ELEMENT NO.	PROJECT NO.	TASK NO.
11. TITLE (Include Security Classification) Are Wait-Free Algorithms Fast?					
12. PERSONAL AUTHOR(S) Hagit Attiya, Nancy Lynch, Nir Shavit					
13a. TYPE OF REPORT Technical		13b. TIME COVERED FROM _____ TO _____		14. DATE OF REPORT (Year, Month, Day) March 1991	
15. PAGE COUNT 40					
16. SUPPLEMENTARY NOTATION					
17. COSATI CODES			18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number)		
FIELD	GROUP	SUB-GROUP			
19. ABSTRACT (Continue on reverse if necessary and identify by block number) The time complexity of wait-free algorithms in "normal" executions, where no failures occur and processes operate at approximately the same speed, is considered. A lower bound of $\log n$ on the time complexity of any wait-free algorithm that achieves <i>approximate agreement</i> among n processes is proved. In contrast, there exists a non-wait-free algorithm that solves this problem in constant time. This implies an $\Omega(\log n)$ time separation between the wait-free and non-wait-free computation models. On the positive side, we present an $O(\log n)$ time wait-free approximate agreement algorithm; the complexity of this algorithm is within a small constant of the lower bound.					
20. DISTRIBUTION/AVAILABILITY OF ABSTRACT <input checked="" type="checkbox"/> UNCLASSIFIED/UNLIMITED <input type="checkbox"/> SAME AS RPT. <input type="checkbox"/> DTIC USERS			21. ABSTRACT SECURITY CLASSIFICATION Unclassified		
22a. NAME OF RESPONSIBLE INDIVIDUAL Carol Nicolora			22b. TELEPHONE (Include Area Code) (617) 253-5894		22c. OFFICE SYMBOL

Are Wait-Free Algorithms Fast?

Hagit Attiya[†]

Nancy Lynch[‡]

Nir Shavit[§]

February 7, 1991

*A preliminary version of this work appeared in the *Proceedings of the 31st Annual Symposium on Foundations of Computer Science, St. Louis*, October 1990. This work was supported by ONR contract N00014-85-K-0168, by NSF grants CCR-8611442 and CCR-8915206, and by DARPA contracts N00014-89-J-1988 and N00014-87-K-0825.

[†]Dept. of Computer Science, Technion, Haifa 32000, Israel. This work was performed while the author was at MIT.

[‡]Laboratory for Computer Science, MIT, Cambridge, MA 02139.

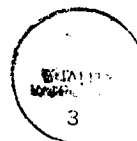
[§]Laboratory for Computer Science, MIT, Cambridge, MA 02139. Part of this work was performed while the author was at the Hebrew University and at the IBM Almaden Research Center.

Keywords: Asynchronous distributed systems, shared memory, wait-free algorithms, read/write atomic registers, lower bounds.

Abstract

The time complexity of wait-free algorithms in "normal" executions, where no failures occur and processes operate at approximately the same speed, is considered. A lower bound of $\log n$ on the time complexity of any wait-free algorithm that achieves *approximate agreement* among n processes is proved. In contrast, there exists a non-wait-free algorithm that solves this problem in constant time. This implies an $\Omega(\log n)$ time separation between the wait-free and non-wait-free computation models. On the positive side, we present an $O(\log n)$ time wait-free approximate agreement algorithm; the complexity of this algorithm is within a small constant of the lower bound.

Accession For	
NTIS CRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A1	



1 Introduction

In shared-memory distributed systems, some number n of independent asynchronous processes communicate by reading and writing to shared memory. In such a computing environment, it is possible for processes to operate at very different speeds, e.g., because of implementation issues such as communication and memory latency, priority-based time-sharing of processors, cache misses and page faults. It is also possible for processes to fail entirely. *Wait-free* algorithms have been proposed as a mechanism for computing in the face of variable speeds and failures: a wait-free algorithm guarantees that each nonfaulty process terminates regardless of the speed and failure of other processes ([23, 28]).¹ The design of wait-free algorithms has been a very active area of research recently (see, e.g., [1, 2, 4, 14, 23, 28, 29, 32, 42, 43, 45, 48]).

Because wait-free algorithms guarantee that fast processes terminate without waiting for slow processes, wait-free algorithms seem to be generally thought of as *fast*. However, while it is obvious from the definition that wait-free algorithms are highly resilient to failures, we believe that the assumption that such algorithms are fast requires more careful examination.

We study the *time complexity* of wait-free and non-wait-free algorithms in “normal” executions, where no failures occur and processes operate at approximately the same speed. We select this particular subset of the executions for making the comparison, because it is only reasonable to compare the behavior of the algorithms in cases where both are required to terminate. Since wait-free algorithms terminate even when some processes fail, while non-wait-free algorithms may fail to terminate in this case, the comparison should only be made in executions in which no process fails, i.e., in *failure-free* executions. The time measure we use is the one introduced in [26, 27], and used to evaluate the time complexity of asynchronous algorithms, in, e.g., [3, 12, 34, 35, 44]. To summarize, we are interested in measuring the time cost imposed by the wait-free property, as measured in terms of extra computation time in the most normal (failure-free) case.

In this paper, we address the general question by considering a specific problem—the *approximate agreement* problem studied, for example, in [15, 19, 20, 36]; we study this problem in the context of a particular shared-memory primitive—single-writer multi-reader atomic registers. In this problem, each process starts with a real-valued input, and (provided it does not fail) must eventually produce a real-valued output. The outputs must all be within a given distance ε of each other, and must be included within the range of the inputs. This problem, a weaker variant of the well-studied problem of distributed consensus (e.g., [21, 30]), is closely related to the important problem of synchronizing local clocks in a distributed system.

Approximate agreement can be achieved very easily if waiting is allowed, by having a designated process write its input to the shared memory; all other processes wait for this value to be written and adopt it as their outputs. In terms of the time measure described above, it is easy to see that the time complexity of this algorithm is constant—independent

¹Wait-free is the shared-memory analogue of the *non-blocking* property for *synchronous* transaction systems (cf. [10, 47]).

of n , the range of inputs and ε . On the other hand, there is a relatively simple wait-free algorithm for this problem, which we describe in Section 3, and which is based on successive averaging of intermediate values. The time complexity of this algorithm depends linearly on n , and logarithmically on the size of the range of input values and on $1/\varepsilon$. A natural question to ask is whether the time complexity of this algorithm is optimal for wait-free approximate agreement algorithms.

Our first major result is an algorithm for the special case where $n = 2$, whose time complexity is constant, i.e., it does *not* depend on the range of inputs or on ε (Section 5). The algorithm uses a novel method of overcoming the uncertainty that is inherent in an asynchronous environment, without resorting to synchronization points (cf. [22]) or other waiting mechanisms (cf. [12]): this method involves ensuring that the two processes base their decisions on information that is approximately, but not exactly, the same.

Next, using a powerful technique of integrating wait-free (but slow) and non-wait-free (but fast) algorithms, together with an $O(\log n)$ wait-free input collection function, we generalize the key ideas of the 2-process algorithm to obtain our second major result: a wait-free algorithm for approximate agreement whose time complexity is $O(\log n)$ (Section 6). Thus, the time complexity of this algorithm does not depend on either the size of the range of input values or on ε , but it still depends on n , the number of processes.

At this point, it is natural to ask whether the logarithmic dependence on n is inherent for wait-free approximate agreement algorithms, or whether, on the other hand, there is a constant-time wait-free algorithm (independent of n). Our third major result shows that the $\log n$ dependency is inherent: any wait-free algorithm for approximate agreement has time complexity at least $\log n$ (Section 7).² This implies an $\Omega(\log n)$ time separation between the non-wait-free and wait-free computation models.

We note that the constant-time 2-process algorithm behaves rather badly if one of the processes fails. The *work* performed in an execution of an algorithm is the total number of atomic operations performed in that execution by all processes before they decide. We present a tradeoff between the time complexity of and the work performed by any wait-free approximate agreement algorithm. We show that for *any* wait-free approximate agreement algorithm for 2 processes, there exists an execution in which the work exhibits a nontrivial dependency on ε and the range of inputs.

In practice, the design of distributed systems is often geared towards optimizing the time complexity in "normal executions," i.e., executions where no failures occur and processes run at approximately the same pace, while building in safety provisions to protect against failures (cf. [31]). Our results indicate that, in the asynchronous shared-memory setting, there are problems for which building in such safety provisions *must* result in performance degradation in the normal executions. This situation contrasts with that occurring, for example, in synchronous systems that solve the distributed consensus problem. In that setting, there are *early-stopping* algorithms (e.g., [16, 18, 40]) that tolerate failures, yet still terminate in *constant* time when no

²The lower bound is attained in an execution where processes run synchronously and no process fails.

failures occur. The exact cost imposed by fault-tolerance on normal executions was studied, for example, in [9, 18, 40]. For synchronous message-passing systems, it has been shown that non-blocking protocols take twice as much time, in failure-free executions, as blocking protocols ([10]).

Recent work has addressed the issue of adapting the usual synchronous shared-memory PRAM model to better reflect implementation issues, by reducing synchrony ([12, 13, 22, 41, 37]) or by requiring fault-tolerance ([25, 24]). To the best of our knowledge, the impact of the *combination* of asynchrony and fault-tolerance (as exemplified by the wait-free model) on the time complexity of shared-memory algorithms has not previously been studied. In [38], Martel, Subramonian and Park present efficient fault-tolerant asynchronous PRAM algorithms. Their algorithms optimize work rather than time and employ randomization. Another major difference is that they assume that inputs are stored in the shared memory, so that every process can access the input of every other process.

The rest of the paper is organized as follows. In Section 2 we present formal definitions of the systems considered in this paper and introduce the time measure. The approximate agreement problem is defined in Section 3, where we also present a fast non-wait-free algorithm and a slow wait-free algorithm for reaching approximate agreement. Section 4 introduces a “bias”-function on which the algorithms in the following sections are based. Proofs of the various properties of this function are, to ease the presentation, deferred to Section 9. A constant time wait-free algorithm for approximate agreement between two processes is presented and proven correct in Section 5; key ideas from this algorithm are used in the $O(\log n)$ time wait-free approximate agreement algorithm presented in Section 6. Section 7 contains the $\log n$ time lower bound for wait-free approximate agreement algorithms. Section 8 presents the lower bound for the tradeoff between the time complexity and the work complexity of a wait-free algorithm for approximate agreement. We conclude, in Section 10, with a discussion of the results and directions for future research.

2 Model of Computation and Time Measure

In this section we describe the systems and the time measure we will consider. Our definitions are standard and are similar to the ones in, e.g., [3, 23, 28, 33, 34].

A *system* consists of n processes p_0, \dots, p_{n-1} . Each process is a deterministic state machine, with a possibly infinite number of states. We associate with each process a set of *local states*. Among the states of each process are a subset called the *initial states* and another subset called the *decision states*. Processes communicate by reading and writing to *single-writer multi-reader atomic registers* R_1, R_2, \dots (also called *shared variables*). Each process p_i has two atomic operations available to it that operate on a shared register R :

- $write(R, v)$ writes the value v to the shared variable R .
- $read(R)$ reads the shared variable R and returns its value v .

A system configuration consists of the states of the processes and the registers. Formally, a *configuration* C is a vector $\langle s_0, \dots, s_{n-1}, v_1, \dots \rangle$ where s_i is the local state of process p_i and v_j is the value of the shared variable R_j . Each shared variable may attain values from some domain which includes a special "undefined" value, \perp . An *initial configuration* is a configuration in which every local state is an initial state and all shared variables are set to \perp . For a configuration $C = \langle s_0, \dots, s_{n-1}, v_1, \dots \rangle$, $state(p_i, C)$ denotes the state of p_i in C and $val(R_j, C)$ denotes the value of register R_j in C , i.e., $state(p_i, C) = s_i$ and $val(R_j, C) = v_j$.

We consider an interleaving model of concurrency, where executions are modeled as sequences of steps. Each step is performed by a single process. A process p_i performs either a *write*(R, v) operation or a *read*(R) operation (which returns a value v), but not both, performs some local computation, and changes to its next local state. The next configuration is the result of these modifications. We assume that each process p_i follows a *local algorithm* A_i that deterministically determines p_i 's next step: A_i determines a variable R and whether p_i is to read or write R as a function of p_i 's local state. If p_i is to read R , then A_i determines p_i 's next state as a function of p_i 's current state and the value v read from R . If p_i is to write R , then A_i determines p_i 's next state and the value v to be written to R as a function of p_i 's current state. An *algorithm* is a function A mapping each i to a local algorithm A_i for p_i .

An *event* on p_i is simply p_i 's index i . A *schedule* is a finite or infinite sequence of events. We denote by λ the empty schedule, with no events. We denote the configuration resulting from the application of a finite schedule σ to a configuration C by $C\sigma$. An *execution fragment* starting from a configuration C is a finite or infinite alternating sequence of configurations and events, $C_0, i_1, C_1, \dots, C_{k-1}, i_k, \dots$, where $C = C_0$ and $C_k = C_{k-1}i_k$, for all $k \geq 1$. We assume that a finite execution fragment ends with a configuration. The *schedule associated with this execution fragment* is i_1, \dots, i_k, \dots . Conversely, the (unique) execution fragment resulting from applying a schedule σ to a configuration C is denoted by (C, σ) . An *execution* is an execution fragment starting with an initial configuration.

Given an infinite schedule σ , a process is *faulty* in σ if it takes a finite number of steps (i.e., has a finite number of events) in σ , and *nonfaulty* otherwise. An infinite schedule σ is *f-admissible* if at most f processes are faulty in σ . In particular, a 0-admissible schedule is called *failure-free*. These definitions also apply to execution fragments by means of their associated schedules.

Let \mathcal{I} be a fixed *input domain* and \mathcal{D} be a fixed *decision domain*. Each initial state of p_i is associated with an input value in \mathcal{I} . For each process p_i and $d \in \mathcal{D}$ we define a subset, $D_{i,d}$, of the states of p_i . We assume that for each p_i , the sets $D_{i,d}$ are pairwise disjoint. We also assume that decisions are irrevocable, i.e., the algorithm transitions are such that if p_i is in a state of $D_{i,d}$ it will remain in a state of $D_{i,d}$. We call the set $D_{i,d}$ the *d-decision states* of p_i .

A *decision problem* (or just *problem*) Π of size n , is a relation between \mathcal{I}^n and \mathcal{D}^n . An algorithm *f-solves* a decision problem Π if in all executions the decisions made can be completed to a decision vector that is in the relation Π to the inputs of the processes. Furthermore, in any *f-admissible* execution, every nonfaulty process eventually decides. An algorithm that

$(n - 1)$ -solves a problem Π is also called a *wait-free* algorithm for Π . Intuitively, even if all processes but one fail when a wait-free algorithm is executed, this process eventually decides.

We now define how to measure the *time* an execution takes.³ We assign times to events in a schedule subject to the following constraints: (a) the time assigned to the first event of any process is at most 1, and (b) the time between two events of the same process is at most 1. The time of a finite schedule σ is the largest amount of real time that can be assigned to the last event in the schedule; denote this by $\text{time}(\sigma)$. The time between two events in a schedule is the largest amount of real time that can elapse between these two events under any time assignment to this schedule. We define the time taken by an execution to be the time taken by the associated schedule. (This definition follows [34, 44].)

An equivalent definition (cf. [3]) is obtained by externally partitioning the computation into minimal rounds: a *round* is any sequence of events such that every process takes a step at least once in the sequence. A *minimal round* is a round such that no proper prefix of it is a round. Every sequence of events can be uniquely partitioned into minimal rounds.⁴ The *time* for an execution is defined to be the number of segments in the unique partition into minimal rounds. (This is the definition introduced in [26, 27], called the *round complexity* in [12].)

The *running time* for p_i in an execution of an algorithm A is defined to be the time associated with the shortest finite prefix of this execution in which p_i is in a decision state (∞ , if there is no such prefix). The *time complexity* of an algorithm A is the supremum of the running times over all failure-free executions of A and all processes p_i .

We conclude this section with some useful notation. Let X be a set of real numbers. Define $\text{range}(X)$ to be the interval $[\min_{x \in X} x, \max_{x \in X} x]$, if X is nonempty and \emptyset , otherwise. Define $\text{diam}(X)$ to be $\max_{x_1, x_2 \in X} |x_1 - x_2|$, if X is nonempty and 0, otherwise. Note that if X is nonempty then $\text{diam}(X)$ is the length of the interval $\text{range}(X)$. If X is nonempty, then $\text{mid}(X) = \frac{\min_{x \in X} x + \max_{x \in X} x}{2}$.

3 Basic Solutions to the Approximate Agreement Problem

We start by defining the *approximate agreement* problem and describing non-wait-free and wait-free algorithms to solve it. In the approximate agreement problem, processes start with real-valued inputs, x_0, \dots, x_{n-1} , and a constant $\epsilon > 0$ (the same ϵ for all processes); all nonfaulty processes are required to decide on real-valued outputs y_0, \dots, y_{n-1} , such that the following conditions hold:

Agreement: for any i, j , $|y_i - y_j| \leq \epsilon$, and

Validity: for any i , $y_i \in \text{range}(\{x_0, \dots, x_{n-1}\})$.

³These definitions can also be formalized in the timed automaton model ([39, 6])

⁴Except, possibly, for the last segment.

<pre> function wait-approx(x); begin 1: $V_0 := x$; 2: return x; end; Process p_0 </pre>	<pre> function wait-approx(x); begin 1: repeat until $V_0 \neq \perp$; /* wait */ 2: return V_0; end; Process $p_i, i \neq 0$ </pre>
--	---

Figure 1: Fast non-wait-free n -process approximate agreement.

This problem has a simple $O(1)$ time non-wait-free solution, described in Figure 1. Process p_0 maintains a single-writer multi-reader atomic register, V_0 , to which it writes its input value as soon as it starts the algorithm. All processes wait until V_0 is set to a value that is not \perp and decide on this value. In the code, any assignment to a shared variable implies a write, and a reference to the value of a shared variable implies a read. Upper case variables denote shared variables, while all lower case variables are local. In this algorithm, the values returned in the return statements are the decision values. Later in the paper, we will use this algorithm as a “subroutine” in our main algorithm; then the values returned in the return statements will not be the final decision values. Similar conventions hold for later algorithms in the paper. We have:

Theorem 3.1 *Procedure wait-approx is a non-wait-free algorithm for the approximate agreement problem whose running time is $O(1)$.*

We next present a wait-free algorithm for approximate agreement. In addition to demonstrating that a wait-free solution exists for this problem, this algorithm will also be used as a “building block” in the construction of a more efficient algorithm, in Section 6.

Let us begin by outlining a simple variant of the algorithm for the case of two processes. Each of the processes $p_i, i \in \{0, 1\}$ has a register which it can write and the other can read. Here and elsewhere, we let \bar{i} denote the index of the other process, i.e., $\bar{i} = 1 - i$. Due to the asynchrony in the system, it is impossible to have processes agree on one of the input values (see [17, 21, 33]). Thus, our algorithm has them gradually converge from the input values x_0 and x_1 to values that are only ε apart. A process p_i repeatedly does the following: It writes its value v_i (initially the input value x_i) into its register, and then reads $p_{\bar{i}}$ ’s register. If p_i reads \perp from $v_{\bar{i}}$, it must decide on its own value, since it can never know when $p_{\bar{i}}$ will write its input value (if at all, because $p_{\bar{i}}$ could have failed before writing). If p_i reads a non- \perp value from $v_{\bar{i}}$, it checks whether or not $|v_{\bar{i}} - v_i| \leq \varepsilon$. If it is, p_i decides on its own value. If not, p_i sets v_i to be $\frac{v_i + v_{\bar{i}}}{2}$ and repeats.

Due to asynchrony, processes do not necessarily converge “directly” to a value. Rather, the following type of scenario is possible: $p_{\bar{i}}$, having formerly written $v_{\bar{i}}$, reads p_i ’s current value v_i , and is delayed just before writing $\frac{v_i + v_{\bar{i}}}{2}$ to its register; then p_i repeatedly reads and writes, cutting the interval in half till its value is very close to $v_{\bar{i}}$; finally, $p_{\bar{i}}$ completes the write of

$\frac{v_i + v_i}{2}$ to its register, so that in fact, p_i moved "too far" towards p_i 's old value. This can repeat itself again and again. However, in every such step of $O(1)$ time (in which both p_i and p_i perform a read and a write), the diameter of the proposed values, $|v_i - v_i|$, is cut by at least a half, and so the values converge in $O(\log(\frac{x_i - x_i}{\epsilon}))$ time. The algorithm is wait-free, since each process can reach a decision independently of the other taking steps.

The algorithm for $n > 2$ processes is of the same flavor, but uses more complicated mechanisms to synchronize among processes. It uses ideas similar to those used in the randomized consensus algorithm of [4]. The computation proceeds in (asynchronous) phases; in each phase, each process suggests a possible decision value. In a manner similar to that of the two process scheme above, the range of suggestions shrinks by a constant factor at each phase, until after $O(\log(\frac{\text{diam}(\{x_0, \dots, x_{n-1}\})}{\epsilon}))$ phases it becomes small enough to allow processes to decide. Because there may be more than two processes, a problem may arise in the case of an execution in which certain slow processes temporarily stop taking steps (i.e. cease advancing in phases), while others (more than one) continue to advance, and then those slow processes return to taking steps again. The algorithm must allow the fast processes to coordinate a decision, while at the same time guaranteeing that the ones that are temporarily slow, will converge to the same decision once they resume activity. The key idea in achieving this task is to allow fast processes that have converged to approximately the same suggested value, and are ahead of all processes with contradictory suggestions by at least two phases, to decide. As will be shown, it can be guaranteed that the processes at lower phases will join this decision value.

The algorithm appears in Figure 2. The inputs to each process p_i are real numbers x_i and ϵ .⁵ For a real number x , define $n_\epsilon(x)$, the ϵ -neighborhood of x , to be $[x - \epsilon, x + \epsilon]$. The algorithm employs a single-writer atomic snapshot object S as a basic memory primitive. Informally, this is a data structure partitioned into n segments S_i , each of which can be updated (written) by one process, and all of which can be scanned (read) by any process in one atomic operation. (More precise specifications and implementations of snapshot objects from single-writer multi-reader atomic registers can be found in [1, 2].) For each process p_i , its segment of S is an array $S_i[1..]$ that in any state contains a finite sequence of reals – its suggestions at different phases – indexed by phase number. Initially, each sequence is λ , the empty sequence. At each phase, after updating (writing) a suggestion to its array (Line 2), a process p_i reads the arrays of all processes (Line 3), obtaining their suggestions for all phases⁶. If p_i is at the maximum phase and all the suggestions by other processes for its phase, or the phase before it, are within ϵ of its latest suggestion, then p_i decides on its latest suggestion (Lines 4-5). Otherwise (Line 6), p_i advances to the next phase taking as its new suggestion the midpoint of all the suggestions at the next phase if there are any, or of its current phase if there are none. Let us make two final remarks before proceeding to prove the algorithm's properties. In the algorithm, since a process first writes to its own sequence and then reads all sequences (including its own), it follows that $\text{phase} \leq \text{max-phase}$. Also, note that in Line 6, r is set to be the number of a phase for which there is at least one suggestion. Thus, the mid operator is applied to a nonempty

⁵Although ϵ is described as a parameter, it is guaranteed that all processes have exactly the same value of ϵ .

⁶Though one can devise algorithms that do not require a process to maintain suggestions for all past phases, we have chosen to do so in order to simplify the exposition and proofs.

```

function wait-free-approx( $x, \epsilon$ );
  begin
1:    $phase := 1$  ;
      repeat forever
2:      $update(S_i[phase] := x)$  ;
3:      $s := scan(S)$  ;
4:      $max-phase := \max_{0 \leq j \leq n-1} \{|s_j|\}$  ; /* Note that  $phase \leq max-phase$  */
5:     if  $phase = max-phase$  and  $phase \geq 2$ 
          and  $s_j[r] \in n_\epsilon[x]$  for all  $j \in \{0, \dots, n-1\}$  and all  $r \geq phase - 1$ 
        then return  $x$  ;
6:     else  $r := \min\{phase + 1, max-phase\}$  ;
7:        $x := mid(\{s_j[r] : |s_j| \geq r\})$  ; /* Note that this set is not empty */
8:        $phase := phase + 1$  ;
    fi;
  end repeat
end;

```

Figure 2: Slow wait-free n -process approximate agreement—Code for process i .

set in Line 7.

We now present the correctness proof for this algorithm. Since the only shared data structure used by the algorithm is the atomic snapshot object S , an execution of the algorithm can be viewed as a sequence of primitive atomic operations that are updates and scans of S . Let α be any execution, and let $r \geq 1$ be a phase number.

For any process $j \in \{0, \dots, n-1\}$, define $S_j^\alpha[r]$ to be the value written by p_j to $S_j[r]$ in α (\perp , if there is no such value). Note that this value is uniquely defined. Define $S^\alpha[r]$ to be $\{S_j^\alpha[r] \neq \perp : j \in \{0, \dots, n-1\}\}$. The following is immediate:

Lemma 3.2 *Let α be an execution and α' is a finite prefix of α . Then $S^{\alpha'}[r] \subseteq S^\alpha[r]$, for every $r \geq 1$.*

Throughout the proofs in this paper, a subscript i for a procedure denotes invocation by process p_i ; similarly, a subscript i for a local variable name denotes the copy of this variable at process p_i . A process p_i is said to be in *phase* r if $phase_i = r$. Denote by $scan_i^r$ the scan performed by p_i at phase r , and by $update_i^r(x)$ the update by p_i at phase r . Note that, for $r \geq 2$, the scan performed before writing a suggestion for phase r is denoted $scan^{r-1}$.

For a finite or infinite execution α and $r \geq 1$, denote

$$mids(\alpha, r) = \{mid(S^{\alpha'}[r]) : \alpha' \text{ is a prefix of } \alpha \text{ and } S^{\alpha'}[r] \text{ is nonempty}\},$$

that is, the set of midpoints of all the sets of suggestions for phase r at earlier points of α . The next lemma is the key for proving that the algorithm is wait-free. It will be used later, in Corollary 3.7, to show that the range of suggestions decreases by a constant factor with each phase. Intuitively, it states that any suggestion for phase r must be in the range of the midpoints of all the sets of suggestions for phase $r - 1$ at earlier points in the execution.

Lemma 3.3 *For any finite execution α and phase $r \geq 2$, $\text{range}(S^\alpha[r]) \subseteq \text{range}(\text{mids}(\alpha, r-1))$.*

Proof: By induction on the length of the execution. The basis holds vacuously.

For the inductive step, the interesting case is when α ends with $\text{update}_i^r(x)$, for some i , where $x = S_i^r[r]$. Then scan_i^{r-1} appears in α . Let α' be the shortest prefix of α that includes scan_i^{r-1} . Note that α' is a proper prefix of α .

Let r' be the largest phase number read in scan_i^{r-1} . Since process p_i reads its own sequence, $r' \geq r - 1$. If $r' = r - 1$, then the code implies that x is the midpoint of $S^{\alpha'}[r - 1]$, which suffices. If $r' \geq r$ then, by the code, $x = \text{mid}(S^{\alpha'}[r])$. By the induction hypothesis on α' , $\text{range}(S^{\alpha'}[r]) \subseteq \text{range}(\text{mids}(\alpha', r - 1))$. Thus,

$$x = \text{mid}(S^{\alpha'}[r]) \in \text{range}(S^{\alpha'}[r]) \subseteq \text{range}(\text{mids}(\alpha', r - 1)) \subseteq \text{range}(\text{mids}(\alpha, r - 1)) ,$$

as needed. ■

Since $\text{range}(\text{mids}(\alpha, r - 1)) \subseteq \text{range}(S^\alpha[r - 1])$, we have:

Corollary 3.4 *For any finite execution α and phase $r \geq 2$, $\text{range}(S^\alpha[r]) \subseteq \text{range}(S^\alpha[r - 1])$.*

For the rest of the proof, we fix some infinite execution β of the algorithm. The following lemmas are stated with respect to β . The following is a corollary of Lemma 3.3.

Corollary 3.5 *For any phase $r \geq 2$, $\text{range}(S^\beta[r]) \subseteq \text{range}(\text{mids}(\beta, r - 1))$.*

The next lemma states that the diameter of all the possible midpoints of the suggestions in phase r is at most half the diameter of all the suggestions for phase r .

Lemma 3.6 *For any phase $r \geq 1$, $\text{diam}(\text{mids}(\beta, r)) \leq \frac{1}{2} \text{diam}(S^\beta[r])$.*

Proof: If $\text{mids}(\beta, r)$ is empty then $\text{diam}(\text{mids}(\beta, r)) = 0$ and the claim follows immediately, so assume that $\text{mids}(\beta, r)$ is nonempty. Let α' and α'' be two prefixes of β such that $S^{\alpha'}[r]$ and $S^{\alpha''}[r]$ are nonempty. It suffices to show that $|\text{mid}(S^{\alpha''}[r]) - \text{mid}(S^{\alpha'}[r])| \leq \frac{1}{2} \text{diam}(S^\beta[r])$. Without loss of generality, suppose α'' is a prefix of α' . By Lemma 3.2, $S^{\alpha''}[r] \subseteq S^{\alpha'}[r] \subseteq S^\beta[r]$. Suppose first that $\text{mid}(S^{\alpha'}[r]) \leq \text{mid}(S^{\alpha''}[r])$. Thus, $\text{mid}(S^{\alpha'}[r]) \leq \text{mid}(S^{\alpha''}[r]) \leq \max(S^{\alpha''}[r]) \leq \max(S^{\alpha'}[r])$. Hence

$$|\text{mid}(S^{\alpha''}[r]) - \text{mid}(S^{\alpha'}[r])| \leq \frac{1}{2} \text{diam}(S^{\alpha'}[r]) \leq \frac{1}{2} \text{diam}(S^\beta[r]) ,$$

as needed. A symmetric argument applies if $\text{mid}(S^{\alpha''}[r]) > \text{mid}(S^{\alpha'}[r])$. ■

The following lemma guarantees that suggestions will become closer with each phase; it will be used together with Lemma 3.9 to ensure wait-freedom.

Lemma 3.7 *For any phase $r \geq 2$, $\text{diam}(S^\beta[r]) \leq \frac{1}{2} \text{diam}(S^\beta[r-1])$*

Proof: By Corollary 3.5, $\text{range}(S^\beta[r]) \subseteq \text{range}(\text{mids}(\beta, r-1))$. Thus,

$$\begin{aligned} \text{diam}(S^\beta[r]) &\leq \text{diam}(\text{mids}(\beta, r-1)) \\ &\leq \frac{1}{2} \text{diam}(S^\beta[r-1]) \quad \text{by Lemma 3.6.} \end{aligned}$$

Lemma 3.8 *If some process returns x in phase r and $y \in S^\beta[r]$, then $y \in n_\epsilon(x)$.*

Proof: Assume p_i returns x in phase r , and assume, by way of contradiction, that there exist processes with suggestions for phase r that are not in $n_\epsilon(x)$. Let p_j be one of these processes with the property that scan_j^{r-1} is the earliest among scan^{r-1} of these processes; let α be the shortest prefix of β that includes scan_j^{r-1} . Let $y = S_j^\beta[r]$; by assumption, $y \notin n_\epsilon(x)$.

By the way p_j was chosen, there is no $\text{update}_j^r(y')$, with $y' \notin n_\epsilon(x)$ in α ; thus, $\text{range}(S^\alpha[r]) \subseteq n_\epsilon[x]$. Let r' be the maximum phase number read in scan_j^{r-1} . It must be that $r' \leq r-1$, since otherwise, p_j would have set its suggestion for phase r to be in $n_\epsilon(x)$. Since process p_j reads its own sequence, $r' = r-1$.

The fact that $r' = r-1$ also implies that scan_j^{r-1} precedes $\text{update}_j^r(x)$. Let α' be the shortest prefix of β that includes scan_j^r . Since $\text{update}_j^r(x)$ precedes scan_j^r , it follows that scan_j^{r-1} precedes scan_j^r , i.e., α is a prefix of α' .

Since process p_i returns in phase r , it follows from the code that $\text{range}(S^{\alpha'}[r-1]) \subseteq n_\epsilon[x]$. Since $r-1$ is the maximum phase number read in scan_j^{r-1} , it follows that $y = \text{mid}(S^{\alpha'}[r-1]) \in \text{range}(S^{\alpha'}[r-1])$. However, by Lemma 3.2, $S^{\alpha'}[r-1] \subseteq S^\alpha[r-1]$, and thus $y \in n_\epsilon(x)$, a contradiction. ■

Lemma 3.9 *For any phase $r \geq 1$, if $\text{diam}(S^\beta[r]) \leq \epsilon$, then every nonfaulty process returns no later than phase $r+1$.*

Proof: From the code of the algorithm it follows that every nonfaulty process either returns or reaches phase $r+1$. If $\text{diam}(S^\beta[r]) \leq \epsilon$ it follows from Corollary 3.4 that $\text{diam}(S^\beta[r+1]) \leq \epsilon$.

The proof proceeds by induction on the order in which processes perform scan^{r+1} . For the base case, let p_i be the first process to perform scan^{r+1} . Clearly, p_i has $\text{phase}_i = r+1 = \text{max-phase}$, and by assumption $r+1 \geq 2$. Also, $\text{diam}(S^\beta[r])$ and $\text{diam}(S^\beta[r+1])$ are less than or equal to ϵ , and thus, p_i will pass the test in Line 5 and will return in phase $r+1$. The induction step is similar, and uses the fact that so far no process has advanced beyond phase $r+1$ to show that any process that reaches phase $r+1$ passes the test in Line 5 and returns in phase $r+1$. ■

Thus we have proved:

Theorem 3.10 *Procedure wait-free-approx is a wait-free algorithm for the approximate agreement problem whose running time on input $\langle x_0, \dots, x_{n-1} \rangle$ is at most*

$$O(n^2 \log(\frac{\text{diam}(\{x_0, \dots, x_{n-1}\})}{\epsilon})) .$$

Proof: The validity condition clearly holds, since processes decide only on their suggestions and these are always within the range of the inputs (Corollary 3.4).

To show agreement, assume that r is the minimum phase in which some process returns, and let p_i be a processes that returns x in phase r . By Lemma 3.8, the suggestions of all processes for phase r are in $n_\epsilon(x)$. By Corollary 3.4, the same is true for phase $r + 1$. By Lemma 3.9, all nonfaulty processes return no later than phase $r + 1$, and thus, all nonfaulty processes return either in phase r or in phase $r + 1$. Since processes return only their suggestions, all returned values are in $n_\epsilon(x)$, as needed.

Since the diameter of suggestions decreases by a factor of two with each phase (by Lemma 3.7) it will eventually be smaller than ϵ and, by Lemma 3.9, each process will eventually decide. This guarantees wait-freedom.

To show the time bound, notice that, by Lemma 3.7, after $O(\log(\frac{\text{diam}(\{x_0, \dots, x_{n-1}\})}{\epsilon}))$ phases, processes will have very close suggestions; by Lemma 3.9, all processes will return. The time it takes a process to execute each phase is bounded from above by the number of operations it executes. Using the implementation of atomic snapshots from [1], this is bounded by $O(n^2)$. ■

Since the input range is not bounded and ϵ may be arbitrarily small, the running time of the algorithm as a function of n is actually unbounded. Note that the time complexity in the execution where processes operate synchronously starting with inputs $\langle x_0, \dots, x_{n-1} \rangle$ is $\Omega(n \log(\frac{\text{diam}(\{x_0, \dots, x_{n-1}\})}{\epsilon}))$.⁷

4 The Bias Function

The algorithms in Sections 5 and 6 return a decision value by performing a calculation based on an input value and a corresponding counter for each process. We name the calculated function bias, as the returned decision value is biased towards (i.e. is closer to) the input value associated with the process having the largest corresponding counter. Before presenting the algorithms, we present the function and explain its properties. The proofs of these properties

⁷The discrepancy between this bound and the bound in the theorem is due to the fact that tighter bounds have not been proven for the time to execute operations in the implementation of atomic snapshot objects of [1].

```

function bias( $v^0, v^1, c^0, c^1, \epsilon$ );
  begin
1:   if  $v^0 = v^1 = 0$  then return 0
2:   else if  $c^0 < c^1$  then return  $v^1 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - \min\{c^1\epsilon, |v^1|\})$ 
3:   else return  $v^0 + \frac{v^1 - v^0}{|v^0| + |v^1|}(|v^0| - \min\{c^0\epsilon, |v^0|\})$ 
    fi;
  end;

```

Figure 3: The bias function—Code for process p_i .

are purely arithmetic, involving no arguments about synchronization between processes, and have therefore been deferred to Section 9.

In order to understand the nature of the calculation performed by the bias function, we briefly explain the structure of the algorithms using it. The new algorithms are conceptually based on the following high-level two-process algorithm. A process p_1 (similarly p_0), knowing only its own input value v^1 , will repeatedly take incremental steps of size ϵ , starting at 0 and ending upon reaching the value v^1 , unless it reads that the other process p_0 has also moved. In the former case it decides on v^1 , and in the latter case its decision value is a function of the relative number of incremental steps both processes managed to take before each noticed the other had moved. However, since in either case process p_1 's decision must be guaranteed to be in $\text{range}(\{v^0, v^1\})$, it cannot just be a value in the interval $\text{range}(\{0, v^1\})$. This is the exact purpose of the function bias. It provides a mapping from the processes' incremental walks in the intervals $\text{range}(\{0, v^0\})$ and $\text{range}(\{0, v^1\})$ respectively, to walks of proportional length in the allowed $\text{range}(\{v^0, v^1\})$. The code of bias appears in Figure 3. The function takes as inputs two real number values v^0 and v^1 , two associated counters, c^0 and c^1 (integers denoting the number of incremental steps each process p_0 or p_1 took), and ϵ .

An example of the translation defined by bias is given in Figure 4 for the case $0 < v^0 < v^1$. Assume p_0 traversed a distance of length $c^0 \cdot \epsilon$ away from 0 towards v^0 , and p_1 a distance of length $c^1 \cdot \epsilon$ away from 0 towards v^1 . The bias function maps the respective distances of length $c^0 \cdot \epsilon$ and $c^1 \cdot \epsilon$ (within the interval $[-v^0, v^1]$), into distances of proportional length in the interval $[v^0, v^1]$. The starting point 0 in $[-v^0, v^1]$, is replaced by the point *new-0* in $[v^0, v^1]$. The returned decision value is then the point associated with the larger counter (larger traversed distance).

We now introduce several lemmas that formally outline the properties of the bias function and on which the correctness proofs of the algorithms in the sequel will be based. The first is a rather simple statement, namely, that the returned value of any call to bias is in $\text{range}(\{v^0, v^1\})$.

Lemma 4.1 *Let c^0, c^1 be nonnegative integers, and v^0, v^1, ϵ be real numbers, with $\epsilon > 0$. Then $\text{bias}(v^0, v^1, c^0, c^1, \epsilon) \in \text{range}(\{v^0, v^1\})$.*

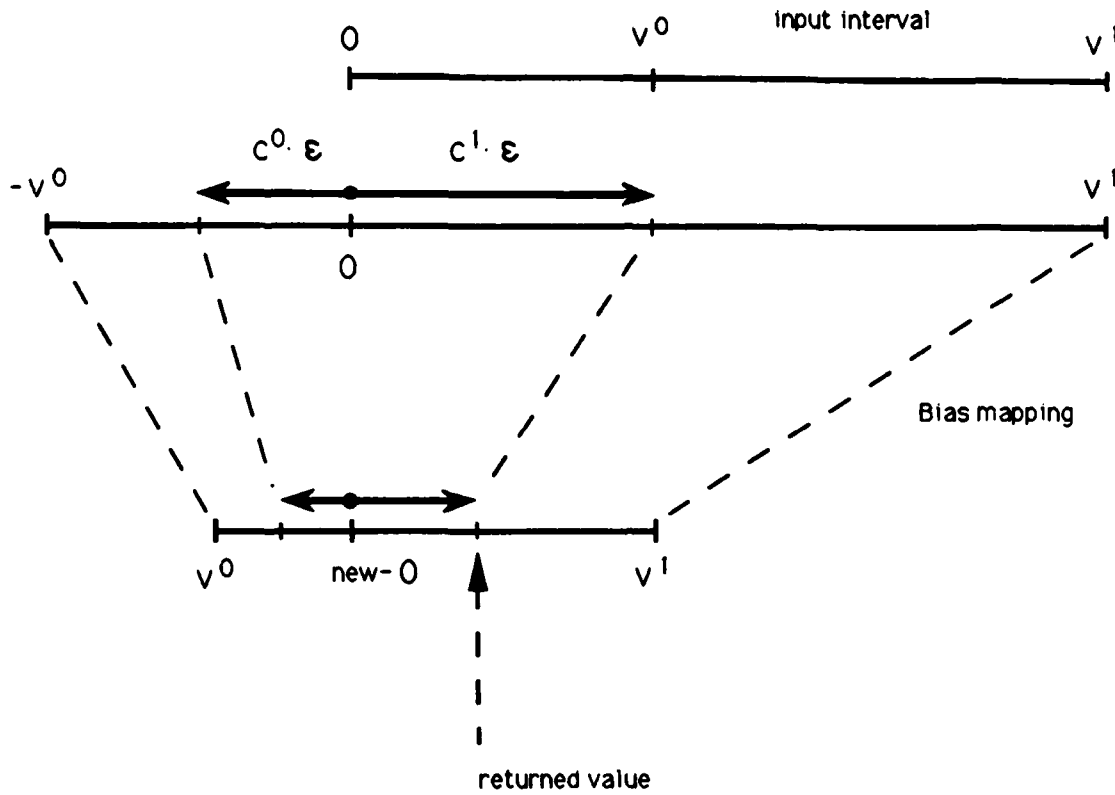


Figure 4: The bias mapping.

The next three lemmas have to do with an additional property required of the bias function: that the values returned by different calls to *bias* always be approximately the same, even if the counter parameter values or the real parameter values used in these calls, are slightly different. The following first lemma states that applying *bias* to counters c^0 and c^1 that are only approximately the same, yet with exactly the same real numbers v^0 , v^1 and ϵ , results in returned values that are approximately the same.

Lemma 4.2 *Let c^0, c^1 be nonnegative integers, and v^0, v^1, ϵ, m be real numbers, $\epsilon > 0$, $m \geq 0$.*

(1) *Suppose $c^1 > c^0$ and $|v^1|/\epsilon - m \leq c^1$. Then $|\text{bias}(v^0, v^1, c^0, c^1, \epsilon) - v^1| \leq m\epsilon$.*

(2) *Suppose $c^0 \geq c^1$ and $|v^0|/\epsilon - m \leq c^0$. Then $|\text{bias}(v^0, v^1, c^0, c^1, \epsilon) - v^0| \leq m\epsilon$.*

The next lemma shows that the results of two calls to *bias* with “close” (in a sense made precise by the lemma) values for c^0, c^1 , and the same v^0, v^1, ϵ , are “close”.

Lemma 4.3 *Let $c_0^0, c_0^1, c_1^0, c_1^1$ be nonnegative integers, and v^0, v^1, ϵ, m be real numbers, $\epsilon > 0$ and $m \geq 0$. Suppose $\min\{c_0^0, c_0^1\} = \min\{c_1^0, c_1^1\} = 0$ and $|c_0^0 - c_0^1| + |c_1^0 - c_1^1| \leq m$. Then*

$$|\text{bias}(v^0, v^1, c_0^0, c_1^0, \epsilon) - \text{bias}(v^0, v^1, c_0^1, c_1^1, \epsilon)| \leq m\epsilon.$$

The last lemma in this section states that applying *bias*, this time to real numbers v^0 and v^1 that are approximately (to within a δ factor) the same, yet with exactly the same counters c^0, c^1 and ϵ , results in values that are approximately the same.

Lemma 4.4 *Let c^0, c^1 be nonnegative integers, and $v_0^0, v_0^1, v_1^0, v_1^1, \varepsilon, \delta$ be real numbers, with $\varepsilon > 0, \delta \geq 0$. Suppose $|v_0^0 - v_1^0| \leq \delta$ and $|v_0^1 - v_1^1| \leq \delta$. Then*

$$|\text{bias}(v_0^0, v_0^1, c^0, c^1, \varepsilon) - \text{bias}(v_1^0, v_1^1, c^0, c^1, \varepsilon)| \leq 6\delta.$$

5 Fast 2-Process Approximate Agreement

We now show that, for two processes, there exists an approximate agreement algorithm whose time complexity is constant; i.e., it does *not* depend on the range of input values or ε . The n -process algorithm presented in Section 6, when specialized to the case $n = 2$, also yields a (somewhat larger) constant time complexity. We present this algorithm because we believe its simplicity will help the reader develop an intuition for the ideas that will be later used in the general algorithm.

5.1 Informal Description

The key ideas underlying this algorithm are as follows. A process, p_i , running on its own, can assume that either it is running very fast (and not much time has elapsed), or the other process, p_r , has failed. Thus, p_i may take an unlimited number of steps without degrading the time complexity for failure-free executions, as long as p_r does not perform any steps. Of course, if p_r does not take any steps at all, then, in order to guarantee the wait-free property, p_i must eventually decide (unilaterally) on its own value. In this case, in order to guarantee correctness, it is necessary that if and when p_r does appear, it must be able to know, just by reading p_i 's registers, what p_i has decided. However, an inherent difficulty of programming asynchronous systems is that, due to the uncertainty of interleaving, at least one process p_i has an "uncertainty of one step," namely, it cannot tell whether p_r read the value written in p_i 's latest write or the value written in p_i 's preceding write. A two-process solution that halves the distance between the suggested values is thus of no use, since the "uncertainty of one step" can cause processes to decide on values that are more than ε apart. Our solution is to have a process change its suggestions gradually with each step, more precisely, by an amount less than ε , so that the "uncertainty of one step" will result only in ε inaccuracy in the decision value.

5.2 The Algorithm

The code for process p_i is given in Figure 5. Each process $p_i, i \in \{0, 1\}$ maintains a *single-writer multi-reader atomic register* with two fields: V_i —the input value, a real number, and C_i —the counter, an integer. Each process starts by writing its input and initializing a counter in the shared memory (Line 1 in `increase-counter`). It then keeps incrementing this counter until either it has taken a number of steps proportional to the absolute value of its input, or the other process has taken a step, whichever happens first (Line 2 of `increase-counter`). When

```

function fast-2-approx( $x, \varepsilon$ );
1:      increase-counter( $x, \frac{\varepsilon}{2}$ ) ;
2:       $\langle v^0, v^1, c^0, c^1 \rangle := \langle V_0, V_1, C_0, C_1 \rangle$ ;
3:      if  $c^i = \perp$  then return  $v^i$ 
4:      else return bias( $v^0, v^1, c^0, c^1, \varepsilon$ );
      end;

function increase-counter( $v, max$ );
1:       $\langle V_i, C_i \rangle := \langle v, 0 \rangle$  ;
2:      while  $C_i = \perp$  and  $C_i < max$  do  $C_i := C_i + 1$  od;
      end;

```

Figure 5: Fast wait-free 2-process approximate agreement—Code for process p_i .

the process stops, it collects all the C and V values and applies the function *bias* to get a decision value. As described in the former section, the decision is within the input range and biased towards the input value of the process with the larger counter. In particular, if a process runs to completion without observing the other process, it decides on its own input value. We show that the discrepancy in the reading of the counters among the two processes is at most 1, and thus, based on the properties of the *bias* function, the decisions based on the values of the counter will differ by at most ε .

5.3 Correctness Proof

An execution of the algorithm can be viewed as a sequence of primitive atomic operations that are reads and writes of atomic registers (and may include changing local data). Fix some execution α of the algorithm. All lemmas in the rest of this section are stated with respect to α . The next lemma shows a crucial property of the values of the counters used by the two processes. In this lemma \perp is treated as -1 .

Lemma 5.1 *Assume p_0 and p_1 return from fast-2-approx. Let $i \in \{0, 1\}$, and let c_i and c_r be the values of C_i read by p_i and p_r , respectively, in Line 2 of fast-2-approx. Then, $c_i - 1 \leq c_r \leq c_i$.*

Proof: Since p_i returns, it must be that p_i writes to C_i . Let π_i be the last write by p_i to C_i in α . Since *increase-counter* returns after the last write to C_i and p_i is the only one to modify C_i , it follows that c_i is the value written to C_i in π_i . Let ϕ_r be the read by p_r of C_i in Line 2 of fast-2-approx. Note that c_r is the value returned in ϕ_r . Since C_i is atomic, it is clear that $c_r \leq c_i$. We now show that $c_i - 1 \leq c_r$.

If $c_i = 0$ then since $c_r \leq c_i$, $c_r \in \{\perp, 0\}$; since \perp is mapped to -1 , the claim follows. So assume $c_i > 0$. Let π'_i be the penultimate write by p_i to C_i , writing $c_i - 1$. Let ϕ_i be the latest

read of C_i by p_i that precedes π_i ; note that π'_i precedes ϕ_i . It must be that the value read in ϕ_i is \perp . Let π_i be the write of 0 by p_i to C_i in α . From the code, it follows that π_i precedes ϕ_i . Since the value read in ϕ_i is \perp , it follows from the atomicity of C_i , that ϕ_i precedes π_i . Thus, π'_i precedes ϕ_i . From the atomicity of C_i it follows that $c_i - 1 \leq c_i$. ■

We can now prove that the algorithm satisfies the agreement property:

Lemma 5.2 *If fast-2-approx₀ returns y_0 and fast-2-approx₁ returns y_1 then $|y_0 - y_1| \leq \varepsilon$.*

Proof: The proof of this lemma is separated into two cases. In one case, we apply Lemma 4.2. In the other case, we show that the sum of the differences between the values of c^0 and c^1 used by p_0 and by p_1 is at most 1, and appeal to Lemma 4.3. The details follow.

Denote by π_i the first write by p_i to C_i , writing 0, for $i \in \{0, 1\}$. Since both processes decide, both π_0 and π_1 must appear in α . Assume, without loss of generality, that π_0 precedes π_1 . (The other case is symmetric.) Assume that process p_0 reads $(v_0^0, v_0^1, c_0^0, c_0^1)$ in Line 2 before deciding, and that process p_1 reads $(v_1^0, v_1^1, c_1^0, c_1^1)$ in Line 2 before deciding. Note that, since p_i first writes 0 to C_i and then reads C_i , it must be that $c_i^i \geq 0$, for $i \in \{0, 1\}$.

Let ϕ be any read of C_0 by p_1 , returning some value z . The code of the algorithm implies that π_1 precedes ϕ . Since π_0 precedes π_1 , π_0 precedes ϕ . By the atomicity of C_0 , this implies that $z \geq 0$. This implies, in particular that $c_1^0 \geq 0$, and thus, fast-2-approx₁ returns in Line 4. In addition, this also implies that p_1 will not increase C_1 beyond 0, and thus, by the atomicity of C_1 , $c_1^1 = 0$ and $c_0^1 \in \{\perp, 0\}$. We separate the rest of the proof into two cases:

Case 1: $c_0^1 = \perp$. In this case, fast-2-approx₀ returns $v_0^0 = x_0$ in Line 3. The code of increase-counter implies that $|x_0|/\varepsilon \leq c_0^0$. Lemma 5.1 implies that $|x_0|/\varepsilon - 1 \leq c_1^0$. Also, $v_1^0 = x_0$. Since $c_1^0 \geq 0 = c_1^1$, we can apply Lemma 4.2(2) with $m = 1$ and get that $|\text{bias}(v_1^0, v_1^1, c_1^0, c_1^1, \varepsilon) - v_0^0| \leq \varepsilon$, as needed.

Case 2: $c_0^1 = 0$. Then fast-2-approx₀ returns in Line 4 and $v_0^1 = v_1^1$. We have that $\min\{c_0^0, c_0^1\} = c_0^1 = 0$ and $\min\{c_1^0, c_1^1\} = c_1^1 = 0$. Also, $|c_0^0 - c_1^0| + |c_0^1 - c_1^1| = |c_0^0 - c_1^0| \leq 1$, by Lemma 5.1. The claim follows by applying Lemma 4.3 with $m = 1$. ■

We have:

Theorem 5.3 *Procedure fast-2-approx is a wait-free algorithm for the 2-process approximate agreement problem whose time complexity is $O(1)$.*

Proof: Agreement follows from Lemma 5.2. It follows from the code and from Lemma 4.1 that the values returned are in the range of the original input values; hence the validity property is satisfied. Each process p_i executes at most $O(|x_i|/\varepsilon)$ steps before deciding; thus, the algorithm is wait-free. Since each process executes a constant number of (its own) steps after the other process performs its first step, the time complexity of this algorithm is $O(1)$. ■

6 Fast n -Process Approximate Agreement

In this section, we present a fast ($O(\log n)$ time) wait-free approximate agreement algorithm for n processes. The algorithm is based on an *alternated-interleaving* method of integrating wait-free (resilient but slow) and non-wait-free (fast but not resilient) algorithms to obtain new algorithms that are both resilient and fast.

We begin by showing how one can reduce, in constant time, the problem of n -process approximate agreement with arbitrary input values to a special case of the problem where the set of input values is included in the union of two small intervals. We do this by performing an alternated-interleaving of a wait-free and a non-wait-free algorithm. We then show, again based on an alternated-interleaving of wait-free and non-wait-free algorithms, that n processes with values in two small intervals can “simulate,” in $O(\log n)$ time, two virtual processes running the fast approximate agreement algorithm of Section 5, thus solving the approximate agreement problem for n processes and any two values. Combining the two algorithms yields an $O(\log n)$ wait-free approximate agreement algorithm.

The second part of the algorithm relies on procedures for synchronization and input collection with $O(\log n)$ time complexity. These procedures are presented in Section 6.3.

6.1 Informal Description

The first part of the algorithm—the one that achieves the constant-time reduction to two small intervals, is encapsulated in procedure *n-to-2* (Figure 6). The idea is simple: interleave the execution of the slow *wait-free-approx* procedure with that of the fast *wait-approx*. The resulting algorithm is wait-free since even if $n - 1$ processes fail, *wait-free-approx* will terminate. It takes at most $O(1)$ time in the failure-free execution since *wait-approx* terminates within $O(1)$ time. However, some processes (group a) might finish the alternated execution with a value from *wait-approx*, while others (group b) finish with a value from *wait-free-approx*. We thus did not solve the approximate agreement problem, but we did guarantee that the values are included in the union of two small intervals. The procedure returns an output value v_i and a group $g_i \in \{a, b\}$ to which p_i is said to belong. It is guaranteed that output values for processes in the same group $g_i \in \{a, b\}$ are at most $\epsilon/12$ apart.

The second part of the algorithm solves n -process approximate agreement in $O(\log n)$ time, assuming that processes are partitioned into two groups with approximately the same value in each group. The solution is based on having the processes in group a (resp. b) jointly simulate a virtual process p_0 (resp. p_1) that execute the function *fast-2-approx* of Figure 5.

The following straightforward simulation is expressed by Lines 1-2 of the function *increase-counter* in Figure 6. The counter C_0 of *fast-2-approx* is replaced by a joint counter, which is defined to be the sum of local counters C_i , for all i in group a . Each step of the simulated counter C_0 is implemented by $O(n)$ steps of the joint counter for a . Each step of this joint counter is, in turn, implemented by a single step of one of the individual counters in group a .

Similarly, the processes in group b simulate counter C_1 of fast-2-approx. In Line 2 of `increase-counter`, in order to decide on the values of the joint counters of a and b , a process reads the values of all local counters. If the counter simulated by p_i 's group is not large enough and the counter simulated by the other group is \perp , then p_i advances the counter simulated by its group (by incrementing its local counter C_i), and repeats. Otherwise, p_i exits `increase-counter`.

One can see that, in an execution where processes operate synchronously, each iteration of the `while` loop in Line 2 of `increase-counter` has $O(n)$ time complexity since reading all memory locations to calculate the simulated counter takes $O(n)$ steps. However, one can improve the time complexity based on the following observation. If p_i ever detects that all processes have set their counters (in Line 1 of `increase-counter`), then it knows that one of the following holds: either some process from the other group has set its local counter (and hence that group's simulated counter), to a value other than \perp , or the other group is empty. In the former case, the loop predicate in Line 2 must be true, while in the latter case, the final value for the other group's counter will be \perp . In either case, p_i can stop executing `increase-counter`, and be guaranteed to correctly simulate the behavior of the 2-process algorithm. In order to detect in less than $O(n)$ time that all processes have set their counters, we use an $O(\log n)$ non-wait-free synch procedure, described in Section 6.3, whose termination ensures this condition. To achieve the better time, the algorithm alternates synch with the (wait-free) loop in Line 2 of `increase-counter`.

The delicate synchronization provided by synch and its effect on the rest of the algorithm guarantee that after some process exits `increase-counter`, individual counter values increase at most by 3. Thus, after exiting `increase-counter`, a process can perform an $O(\log n)$ wait-free fast-collect, described in Section 6.3, in order to collect all the values needed to decide on the returned value in Lines 3-4. The above property ensures that the simulated counter values used by different processes do not differ much.

6.2 The Algorithm

The code for the algorithm is presented in Figure 6. Alternated procedures are enclosed within `begin-alternate` and `end-alternate` brackets. This construct means that the algorithm alternates strictly between executing single steps of the two alternated procedures, and terminates the first time one of the procedures terminates.⁸ When an alternation is used in an assignment statement, the value assigned is the value returned by the procedure that terminates first. The algorithm uses the bias procedure of Figure 3. In addition to the shared data structures used by wait-free-approx and wait-approx, process $p_i, i \in \{0, \dots, n-1\}$, has a *single-writer multi-reader atomic register* with the following fields: V_i —the value returned in p_i 's first phase; G_i —denoting the group to which p_i belongs; C_i — p_i 's contribution to its group's counter; T_i — p_i 's boolean synch termination flag.

⁸We remark that this is just a coding convenience, used to simplify the control structure of the algorithm. It is implemented locally at one process and does not cause spawning of new processes.

In the code we abuse notation and denote by V^g , where g is a group's name, the "group's value" calculated as follows: if $g = g_i$ then it is V_i , and if $g \neq g_i$ then it is an arbitrary V_j such that p_j is in group g if there is any, and \perp , otherwise. The value v^g is calculated in a similar manner from the corresponding local copies. (Recall our convention that lower case letters stand for local variables and upper case letters for shared variables.) When g is a group name, \bar{g} denotes the other group's name, e.g., if $g = a$ then $\bar{g} = b$. The notation C^g , for $g \in \{a, b\}$, stands for the sum of those C_i such that $G_i = g$ and $C_i \neq \perp$, if there is any such C_i , and \perp , otherwise. The value c^g is calculated in a similar manner from the corresponding local copies.

6.3 Fast Information Collection and Synchronization

We now present the procedures for information collection and synchronization and prove their properties. We start with a wait-free algorithm for *input collection*—returning the current values in the entries of an array R . The time complexity of the algorithm is $O(\log n)$.

This problem is interesting on its own as it underlies any problem of computing a function, e.g., max or sum, on a set of initial values that reside in the shared memory.⁹ Once a process collects all the values, computing the function can be done locally in constant time. Since $\Omega(\log n)$ is a lower bound on the time for the information collection problem (see, e.g., [11]), this implies that for problems whose output depends on all the initial values in memory, and only on them, there exists an optimally fast wait-free solution.

Our algorithm, presented in Figure 7, is a wait-free variation of the *pointer-jumping* technique used in PRAM algorithms (e.g., [49]). For sequences R, R' and a nonnegative integer n we define *concatenate* (R, R') as returning the concatenation of R' to R , and *truncate* (R, n) as returning the first n elements of R if $|R| > n$, and R , otherwise. The initial value \perp is treated like any other value and may be returned by the algorithm for entries that have not yet been set.

Fix some execution α of fast- n -approx algorithm. We clearly have:

Theorem 6.1 *Assume fast-collect_i is invoked by p_i in α , and let α' be the shortest prefix of α that includes an invocation of fast-collect. Then fast-collect_i returns a vector containing, for each p_j , a value that appears in R_j at some point at or after α' . Moreover, fast-collect_i returns within at most $2n$ steps by p_i .*

The next lemma is the crux of the time analysis for this algorithm.

Let t be the time of the last event in the shortest finite prefix of α that includes an invocation of fast-collect by every p_i , $i \in \{0, \dots, n-1\}$, if such a prefix exists, ∞ otherwise.

⁹Note that these problems are very different from the *decision problems* considered until now in this paper, where inputs are local to the processes and do not reside in the shared memory.

```

function fast-n-approx( $x, \varepsilon$ );
  begin
0:    $\langle v, g \rangle := \text{n-to-2}(x, \varepsilon)$  ;
1:    $\text{increase-counter}(v, g, \frac{|v|}{\varepsilon/6n})$  ;
2:    $\langle \bar{v}, \bar{g}, \bar{c} \rangle := \text{fast-collect}(V, G, C)$ ;
3:   if  $c^g = \perp$  then return  $v^g$ 
4:   else return  $\text{bias}(v^a, v^b, c^a, c^b, \varepsilon/6n)$ ;
  end;

function n-to-2( $x, \varepsilon$ );
  begin
     $\langle v, g \rangle := \text{begin-alternate}$ 
1:    $\langle \text{wait-free-approx}(x, \varepsilon/12), a \rangle$ 
      and
2:    $\langle \text{wait-approx}(x), b \rangle$ ;
    end-alternate;
3:   return  $\langle v, g \rangle$ 
  end;

function increase-counter( $v, g, \text{max}$ );
  begin
1:    $\langle V_i, G_i, C_i \rangle := \langle v, g, 0 \rangle$ ;
    begin-alternate
2:   while  $C^g = \perp$  and  $C^g < \text{max}$  do  $C_i := C_i + 1$  od;
    and
3:    $\text{synch}(C)$ ;
    end-alternate;
4:    $T_i := \text{true}$ ;
  end;

```

Figure 6: Fast wait-free n -process approximate agreement—Code for process p_i .


```

function fast-collect ( $R$ );
  begin
1:       $l := 1$ ;
2:      while  $l < n$  do
3:           $R_i := \text{concatenate}(R_i, R_{(i+l) \bmod n})$ ;
4:           $l := |R_i|$ ;
        od;
5:      return truncate( $R_i, n$ );
  end;

```

Figure 7: Fast wait-free information collection—Code for process p_i .

Lemma 6.2 Assume $t < \infty$. For every $i \in \{0, \dots, n-1\}$ and every integer r , $0 \leq r \leq \lceil \log n \rceil$, $|R_i| \geq \min\{2^r, n\}$ at time $t + 2r$.

Proof: The proof is by induction on r . The base case, $r = 0$, is trivial.

For the induction step, assume that $r \geq 1$. If at time $t + 2(r-1)$, $|R_i| \geq n$, the claim follows. So suppose, $|R_i| < n$ at time $t + 2(r-1)$. Then process p_i reads $R_{(i+l_i) \bmod n}$ after time $t + 2(r-1)$, where l_i is the length of R_i at time $t + 2(r-1)$. By the inductive hypothesis, $|R_i| \geq 2^{r-1}$ and $|R_{(i+l_i) \bmod n}| \geq 2^{r-1}$, at time $t + 2(r-1)$. It follows that $|R_i| \geq 2^r$ at time $t + 2r$. ■

In particular, at time $t + 2\lceil \log n \rceil$, we have $|R_i| \geq n$ for every i . Thus, fast-collect _{i} returns by time $t + 2\lceil \log n \rceil$. We have:

Theorem 6.3 Let α' be a finite prefix of α . Assume that in α' , fast-collect _{i} is invoked by p_i , for every $i \in \{0, \dots, n-1\}$. Then for every $i \in \{0, \dots, n-1\}$, fast-collect _{i} returns within at most $O(\log n)$ time after α' .

The synchronization procedure, *synch*, is a variant of fast-collect. Since it is used within an *alternate* construct, it is possible that *synch* is aborted without completing and returning “normally.” To cope with this possibility, we associate with the shared array R to which *synch* is applied, a special *termination* array T , whose entries can take on values $\{\perp, \text{true}\}$. T_j is set to *true* if p_j terminates, i.e., aborts or returns from *synch* _{j} . The synchronization procedure guarantees that if it returns, then either all the entries of the array are non- \perp values, or for some j , $T_j = \text{true}$. It is *not* wait-free. The code appears in Figure 8.

Again, we fix some execution α of fast-n-approx. We have:

```

procedure synch( $R$ );
  begin
    1:      repeat until  $R_i \neq \perp$ ;                                /* wait */
    2:       $l := 1$ ;
    3:      while  $l < n$  and  $T_{i+l \bmod n} = \perp$  do
    4:        repeat until  $R_{i+l \bmod n} \neq \perp$ ;                    /* wait */
    5:         $R_i := \text{concatenate}(R_i, R_{(i+l) \bmod n})$ ;
    6:         $l := |R_i|$ ;
      od;
  end;

```

Figure 8: Fast non-wait-free synchronization—Code for process p_i .

Theorem 6.4 *Let α' be a finite prefix of α . Assume that in α' all R entries are set to values $\neq \perp$, and that synch_i is invoked by p_i . Then synch_i returns within at most $3n$ steps by p_i after the end of α' .*

Theorem 6.5 *Let α' be a finite prefix of α . Assume that in α' , synch_i returns, for some p_i . Then at the end of α' either all R entries are $\neq \perp$ or $T_j = \text{true}$ for some j .*

Let α' be a finite prefix of α . Note that, in fast- n -approx, if p_i terminates synch_i , i.e., aborts or returns, then within one time unit, $T_i = \text{true}$. This is crucial in the proof of the next theorem.

Theorem 6.6 *Let α' be a finite prefix of α . Assume that in α' all R entries are set to values $\neq \perp$, and synch_i is invoked by p_i , for every $i \in \{0, \dots, n-1\}$. Then every process terminates synch within at most $O(\log n)$ time after the end of α' .*

Proof: Let t be the time of the last event of α' . We prove that for every process p_i and for every integer r , $0 \leq r \leq 2\lceil \log n \rceil$, by time $t + 3r$, either p_i terminates synch_i or $|R_i| \geq \min\{2^{\lfloor r/2 \rfloor}, n\}$. The claim follows by taking $r = 2\lceil \log n \rceil$ and noticing that if $|R_i| \geq n$, then p_i returns from synch_i within $O(1)$ time.

The proof is by induction on r . The base cases, $r = 0, 1$, are trivial.

For the induction step, assume that $1 \leq r \leq 2\lceil \log n \rceil$. If p_i terminates by time $t + 3r$, then the claim is immediate. So, assume p_i does not terminate by time $t + 3r$. In particular, it does not terminate by time $t + 3(r-1)$. Hence, by the induction hypothesis, $|R_i| \geq \min\{2^{\lfloor (r-1)/2 \rfloor}, n\} = 2^{\lfloor (r-1)/2 \rfloor}$. Then process p_i reads $T_{(i+l_i) \bmod n}$ after time $t + 3(r-1)$, where l_i is the length of R_i at time $t + 3(r-1)$.

If $p_{(i+l_i) \bmod n}$ terminates by time $t + 3(r - 1) - 1$, then, by assumption, $T_{(i+l_i) \bmod n} = \text{true}$ by time $t + 3(r - 1)$ and thus, p_i terminates by time $t + 3r$. It follows from the induction hypothesis for $r - 2$ that $|R_{(i+l_i) \bmod n}| \geq 2^{\lfloor (r-2)/2 \rfloor}$. Then the length of R_i at time $t + 3r$ is larger than or equal to $2^{\lfloor (r-1)/2 \rfloor} + 2^{\lfloor (r-2)/2 \rfloor} \geq 2^{1+\lfloor (r-2)/2 \rfloor} = 2^{\lfloor r/2 \rfloor}$. ■

6.4 Correctness Proof

We remind the reader that an execution of the algorithm is viewed as a sequence of primitive atomic operations that are reads and writes of atomic registers. We now fix some execution α of fast-n-approx.

As in the proof of the 2-process algorithm (Section 5), the crucial point in the proof of the algorithm is showing that, in Lines 3-4 of fast-n-approx, processes use "close" values for c^a and c^b . We show that the value of an arbitrary counter when some process invokes fast-collect are at most 3 less than the maximal value this counter ever attains. This is formalized and proved in the next lemma:

Lemma 6.7 *Assume that p_i invokes fast-collect_i in α . Fix some process p_j ; let k be the value of C_j returned by fast-collect_i. Let k' be the maximum value written to C_j in α . Then $k' - 3 \leq k \leq k'$.*

Proof: The inequality $k \leq k'$ follows immediately from the atomicity of the shared register. To prove the other inequality, let $p_{i'}$ be the first process to execute the write operation in Line 4 of increase-counter. Such a process exists because p_i performs this write operation before invoking fast-collect_i. Let α' be a shortest prefix of α that includes $p_{i'}$'s write to $T_{i'}$. Let k'' be the value of C_j at the end of α' . Since any invocation of fast-collect follows this last write operation in Line 4, Theorem 6.1 and the atomicity of C_j implies that $k'' \leq k$. Thus, it suffices to show that $k' - 3 \leq k''$. There are two cases according to the way p_i exits the alternate construct in Lines 2-3 of increase-counter:

Case 1: $p_{i'}$ exits the while loop. It must be that one of the halting conditions of the while loop is false for $p_{i'}$. If $p_{i'}$ and p_j are in the same group, i.e., $g_{i'} = g_j$, then p_j will perform at most one iteration of the while loop before p_j also sees the corresponding condition to be false. If $p_{i'}$ and p_j are not in the same group, i.e., $g_{i'} \neq g_j$, then p_j will perform at most one iteration of the while loop before p_j sees the first condition to be false (by observing $C_{i'} \neq \perp$). The claim follows.

Case 2: $p_{i'}$ returns from synch_i. It follows that for all processes, $T_j = \perp$ when p_i terminates synch_i. It follows from Theorem 6.5 that, for all $l \in \{0, \dots, n - 1\}$, the value of C_l at the end of α' is $\neq \perp$. By Theorem 6.4, p_j will exit synch_j(C) after performing at most $3n$ of its own steps after α' . It follows from the definition of alternate that p_j will perform at most $3n$ steps in the while loop in Line 2 of increase-counter, before synch_j(C) terminates. However, each iteration of the while loop takes at least n steps (since n registers have to be read).

Thus, p_i will perform at most three additional iterations of the **while** loop, before $\text{synch}_j(C)$ terminates. The claim follows. \blacksquare

This implies that, for each local counter, the values read by two different processes differ at most by 3. Hence, the values used by different processes for the joint counters c^a and c^b differ at most by $3n$. Formally, we have:

Lemma 6.8 *Suppose $i, j \in \{0, \dots, n-1\}$ and $g \in \{a, b\}$. Assume the values returned by fast-collect_i and fast-collect_j calculate to c_i^g and c_j^g , respectively. Then $|c_i^g - c_j^g| \leq 3n$.*

We can now prove that the algorithm satisfies the agreement property:

Lemma 6.9 *If fast-approx_i returns y_i and fast-approx_j returns y_j , then $|y_i - y_j| \leq \epsilon$.*

Proof: The general outline of the proof parallels that of Lemma 5.2; however, some of the details are different. First, the discrepancy between processes' view of the joint counters might be $3n$; to compensate for that, we use bias with $\epsilon/6n$. In addition, we must allow for the possibility of using different values from the same group (by applying Lemma 4.4). The details follow.

We present the proof for the case where p_i and p_j are not in the same group, without loss of generality, assume $g_i = a$ and $g_j = b$. The proof for the case where p_i and p_j are in the same group follows from similar arguments and is left to the reader.

Assume that the values computed by p_i based on fast-collect_i to be used in Lines 3-4 of fast-n-approx are $(v_i^a, v_i^b, c_i^a, c_i^b)$; similarly, assume that the values computed by p_j based on fast-collect_j to be used in Lines 3-4 of fast-n-approx are $(v_j^a, v_j^b, c_j^a, c_j^b)$. Note that since p_i is in group a , $c_i^a \geq 0$ and $v_i^a \neq \perp$; similarly, since p_j is in group b , $c_j^b \geq 0$ and $v_j^b \neq \perp$.

For any process p_k , denote by π_k the write by process p_k in Line 1 of increase-counter (if it appears in α). Since p_i and p_j decide, π_i and π_j must appear in α . Let $p_{i'}$ be such that $\pi_{i'}$ is the first π_k in α . Assume, without loss of generality, that $p_{i'}$ is in group a . Intuitively, we assume that the first process to start the second phase of the algorithm belongs to p_i 's group, a .

The code of the algorithm implies that $\pi_{i'}$ precedes any calculation of C^a by $p_{j'}$, for any $p_{j'}$ in group b . Since $\pi_{i'}$ precedes π_j , it follows that $p_{j'}$ will always calculate $C^a \neq \perp$. Thus, $c_j^a \geq 0$ and hence fast-n-approx_j returns in Line 4 and $v_j^a \neq \perp$. Also, the above implies that C^b never increases beyond 0. Thus, $c_j^b = 0$ and $c_i^b \in \{\perp, 0\}$. We separate the rest of the proof into two cases:

Case 1: $c_i^b = \perp$. Then fast-n-approx_i returns v_i^a in Line 3. From the code it follows that $c_i^a \geq |v_i^a|6n/\epsilon$. By Lemma 6.8, $c_j^a \geq |v_i^a|6n/\epsilon - 3n$. Since $c_j^a \geq 0 = c_j^b$, applying Lemma 4.2 (2) with $m = 3n$ we get that

$$|\text{bias}(v_i^a, v_j^b, c_j^a, c_j^b, \epsilon/6n) - v_i^a| \leq \epsilon/2. \quad (1)$$

Also, Theorem 3.1 implies that $|v_i^a - v_j^a| \leq \epsilon/12$. Applying Lemma 4.4 with $\delta = \epsilon/12$, $c^0 = c_j^a$, $c^1 = c_j^b$, $v_0^0 = v_j^a$, $v_0^1 = v_j^b$, $v_1^0 = v_i^a$, $v_1^1 = v_j^b$, we get that

$$|\text{bias}(v_j^a, v_j^b, c_j^a, c_j^b, \epsilon/6n) - \text{bias}(v_i^a, v_j^b, c_j^a, c_j^b, \epsilon/6n)| \leq 6\epsilon/12 = \epsilon/2. \quad (2)$$

From (1) and (2) it follows that

$$|\text{bias}(v_j^a, v_j^b, c_j^a, c_j^b, \epsilon/6n) - v_i^a| \leq \epsilon,$$

as needed.

Case 2: $c_i^b = 0$. Thus, `fast-n-approx` returns in Line 4 and $v_i^b \neq \perp$. We have that $\min\{c_i^a, c_i^b\} = c_i^b = 0$ and $\min\{c_j^a, c_j^b\} = c_j^b = 0$. Also, $|c_i^a - c_j^a| + |c_i^b - c_j^b| = |c_i^a - c_j^a| \leq 3n$ by Lemma 6.8. Applying Lemma 4.3 with $m = 3n$ we get

$$|\text{bias}(v_j^a, v_j^b, c_i^a, c_i^b, \epsilon/6n) - \text{bias}(v_j^a, v_j^b, c_j^a, c_j^b, \epsilon/6n)| \leq 3n \cdot \epsilon/6n = \epsilon/2. \quad (3)$$

Also, Theorems 3.1 and 3.10 imply that $|v_i^a - v_j^a| \leq \epsilon/12$ and $|v_i^b - v_j^b| \leq \epsilon/12$. By applying Lemma 4.4 with $\delta = \epsilon/12$ we get

$$|\text{bias}(v_i^a, v_i^b, c_i^a, c_i^b, \epsilon/6n) - \text{bias}(v_j^a, v_j^b, c_i^a, c_i^b, \epsilon/6n)| \leq 6\epsilon/12 = \epsilon/2. \quad (4)$$

From (3) and (4) it follows that

$$|\text{bias}(v_i^a, v_i^b, c_i^a, c_i^b, \epsilon/6n) - \text{bias}(v_j^a, v_j^b, c_j^a, c_j^b, \epsilon/6n)| \leq \epsilon,$$

as needed. ■

We have:

Theorem 6.10 *Procedure fast-n-approx is a wait-free algorithm for the n-process approximate agreement problem whose time complexity is $O(\log n)$.*

Proof: Agreement follows from Lemma 6.9. Validity follows immediately since the values returned by `wait-free-approx` and `wait-approx` are in the range of the original inputs, and the `bias` function preserves this property (Lemma 4.1).

The algorithm is wait-free because the first alternative of each alternation construct and `fast-collect` are wait-free.

Within $O(1)$ time all processes finish `n-to-2`. Thus, within $O(1)$ time all processes start `procedure increase-counter`, write to C_i and invoke `synch`. By Theorem 6.6, within $O(\log n)$ time each process terminates `synch`. Thus, within $O(\log n)$ time all processes exit `increase-counter` and invoke `fast-collect`. By Theorem 6.3, all processes return from `fast-collect` within $O(\log n)$ time. Hence, the total time complexity is $O(\log n)$. ■

7 A $\log n$ Time Lower Bound

In this section, we show that the $\log n$ dependency exhibited by the algorithm of Theorem 6.10 is inherent: the time complexity of any wait-free algorithm for n -process approximate agreement is at least $\log n$. Together with Theorem 3.1, this result shows that there are problems for which wait-free algorithms take more time (by an $\Omega(\log n)$ factor) than non-wait-free algorithms.

In the rest of this section, we assume that each process has only one register to which it can write. Since the size of registers is not restricted and since only one process may write to each register, there is no loss of generality in this assumption. Let R_i be the register to which p_i writes. For a configuration C and a process p_i , let $st(p_i, C)$ be the pair consisting of the local state of p_i and the value of R_i in C , i.e., $st(p_i, C) = \langle state(p_i, C), val(R_i, C) \rangle$.

The *synchronized schedule* is the schedule in which processes take steps in round-robin order starting with p_0 , essentially operating synchronously. The sequence of r rounds in the round-robin order is denoted σ_r . For any configuration C , the corresponding *synchronized execution* from C is uniquely determined by the algorithm. Note that this is a 0-admissible execution.

We now define the set of processes that could have influenced p_i 's state at time r in the synchronized execution from a configuration C . Let C be a configuration; by induction on $r \geq 0$, define the set $INF(p_i, r, C)$, for every $i \in \{0, \dots, n-1\}$, using the following rules:

1. $r = 0$: $INF(p_i, r, C) = \{p_i\}$, for every $i \in \{0, \dots, n-1\}$.
2. $r \geq 1$: if p_i 's r th step in (C, σ_r) is a read of R_j , then $INF(p_i, r, C) = INF(p_i, r-1, C) \cup INF(p_j, r-1, C)$. If p_i 's r th step is a write (to R_i) then $INF(p_i, r, C) = INF(p_i, r-1, C)$.

The next lemma formalizes the intuition that INF includes all the processes that can influence p_i 's state up to time r .

Lemma 7.1 *Let C_1 and C_2 be two configurations, let p_i be any process and let $r \geq 0$. If $st(p_j, C_1) = st(p_j, C_2)$ for all $p_j \in INF(p_i, r, C_1)$, then $st(p_i, C_1\sigma_r) = st(p_i, C_2\sigma_r)$.*

Proof: The proof is by induction on r . The base case, $r = 0$, is trivial since in this case, $\sigma_0 = \lambda$, $INF(p_i, 0, C_0) = \{p_i\}$ and the claim follows from the assumptions.

To prove the induction step, assume the claim holds for $r-1$. If p_i 's r th step is a write then the claim follows immediately from the induction hypothesis since $INF(p_i, r-1, C_1) = INF(p_i, r, C_1)$.

If p_i 's r th step is a read from R_j then $INF(p_j, r-1, C_1) \subseteq INF(p_i, r, C_1)$. The induction hypothesis implies that $st(p_j, C_1\sigma_{r-1}) = st(p_j, C_2\sigma_{r-1})$. By the same reasoning, $st(p_i, C_1\sigma_{r-1}) = st(p_i, C_2\sigma_{r-1})$. Thus, $st(p_i, C_1\sigma_r) = st(p_i, C_2\sigma_r)$, as needed. ■

We can now prove:

Theorem 7.2 *Any wait-free algorithm for the n -process approximate agreement problem has time complexity at least $\log n$.*

Proof: Assume that A is a wait-free approximate agreement algorithm. We prove a slightly stronger claim: there exists a 0-admissible execution α in which no process decides before time $\log n$. Suppose, by way of contradiction, that in all 0-admissible executions some process decides before time $\log n$.

Fix some $\varepsilon < 1$. Let σ be the infinite synchronized schedule, i.e., the limit of σ_r . Consider the execution of A under σ from the initial configuration C_0 where processes start with inputs $(0, \dots, 0)$. Let t be the time associated with the first decision event in (C_0, σ) ; without loss of generality, let p_0 be the process associated with this event. By assumption, $t < \log n$. By the validity property, p_0 must decide on 0 since all processes start with 0. Note that, by induction on r , $|INF(p_i, r, C)| \leq 2^r$, for every configuration C , $r \geq 0$ and $i \in \{0, \dots, n-1\}$. Since $t < \log n$ it must be that $|INF(p_0, T, C_0)| \leq 2^T < n$. Thus, there exists some process that is not in $INF(p_0, t, C_0)$; without loss of generality, assume $p_{n-1} \notin |INF(p_0, T, C_0)|$.

Intuitively, to complete the proof, we create an alternative execution in which p_{n-1} "starts early" with input 1, runs on its own and thus must eventually decide on 1. We then let the rest of the processes execute as if they are in the synchronized execution from C_0 and use Lemma 7.1 to show that process p_0 still decides on 0, which is a contradiction to the agreement property, since $\varepsilon < 1$.

More precisely, apply τ , an infinite schedule consisting of steps of p_{n-1} only, to the initial configuration C_2 , where processes start with inputs $(1, \dots, 1)$. The resulting execution (C_2, τ) is $(n-1)$ -admissible, and thus, since A $(n-1)$ -solves the approximate agreement problem, and since p_{n-1} is nonfaulty in τ , there exists a finite prefix τ' of τ in which p_{n-1} decides. By validity, p_{n-1} decides on 1. Now apply τ' to the initial configuration C_1 where all processes but p_{n-1} start with input 0, and p_{n-1} starts with input 1. By induction on the prefixes of τ' , it follows that the $st(p_{n-1}, C_1\tau') = st(p_{n-1}, C_2\tau')$. Thus p_{n-1} decides on 1 in $C_1\tau'$. Since p_{n-1} can write only to R_{n-1} , it follows that for all processes $p_j \neq p_{n-1}$, $st(p_j, C_1\tau') = st(p_j, C_0)$. By Lemma 7.1, $state(p_0, C_1\tau'\sigma_T) = state(p_0, C_0\sigma_T)$. Thus, p_0 decides on 0 in $C_1\tau'\sigma_T$, and p_{n-1} decides 1, which is a contradiction to agreement, since $\varepsilon < 1$. ■

8 A Tradeoff Between Work and Time

We now consider the performance of wait-free algorithms when failures occur. A drawback of the fast algorithms we have presented in this paper is that if a failure *does* occur, then the remaining processes will have to take many steps before halting. We show that this phenomenon is unavoidable. Roughly speaking, we prove that if an algorithm terminates

in a small number of steps in executions where failures do occur, then it is slow in normal executions. In the rest of this section we restrict our attention to the 2 processes case.

Recall that the work performed by an algorithm is define to be the maximum, over all executions, of the total number of operations performed by all processes before deciding. The lower bound presented here is slightly stronger—it bounds the number of operations a single process performs before deciding when running on its own. Clearly, this also gives a lower bound on the work.

Let $k \geq 1$ be an integer. An algorithm is k -bounded if from any reachable configuration, a process that executes k consecutive steps on its own must decide. Fix a k -bounded wait-free algorithm A for approximate agreement; all definitions and lemmas in the rest of this section are with respect to A . For each process p_i and a configuration C , reachable in an execution of A , define $\text{pref}_i(C)$, the preference of p_i in C , to be the value on which p_i decides in the execution fragment starting from C in which it runs alone until it decides.

A finite schedule is a *block* if it consists of a positive number of events by p_0 followed by one event by p_1 , or vice versa.

Lemma 8.1 *Let σ be a finite schedule, and let C_0 be an initial configuration. Let $C = C_0\sigma$. There exists a finite block schedule σ' such that*

$$|\text{pref}_0(C\sigma') - \text{pref}_1(C\sigma')| \geq \frac{1}{2k} |\text{pref}_0(C) - \text{pref}_1(C)|.$$

Proof: The proof considers the tree of all execution fragments of time 1 from C . A case analysis, according to the types of steps taken, similar to the one in [33], is used to show that it cannot be that all the pairs of preferences associated with leaves of this tree are close together. The details follow.

Let $\tau_0 = 0^k$, i.e., the schedule consisting of k events of p_0 . Similarly, let $\tau_1 = 1^k$. Let $(C, \tau_0) = C, C_1, \dots, C_k$, and $(C, \tau_1) = C, C'_1, \dots, C'_k$. For any l , $1 \leq l \leq k$, denote $D_l = C_l 1$, i.e., the configuration that results from applying an event of p_1 to C_l . Similarly, for any l , $1 \leq l \leq k$, denote $D'_l = C'_l 1$. Denote $v_0^l = \text{pref}_0(D_l)$, $v_1^l = \text{pref}_1(D_l)$, $u_0^l = \text{pref}_0(D'_l)$ and $u_1^l = \text{pref}_1(D'_l)$.

Since A is k -bounded, it must be that p_0 decides in $C\tau_0$; by definition, it must decide on $\text{pref}_0(C)$. Similarly, p_1 decides on $\text{pref}_1(C)$ in $C\tau_1$. We show that for all l , $1 \leq l < k$, either $v_0^l = v_0^{l+1}$ or $v_1^l = v_1^{l+1}$. There are four cases, depending on the type of operation taken in p_0 's step from C_l to C_{l+1} and in p_1 's step from C_l to D_l :

1. p_0 writes and p_1 writes: commutativity implies that $v_0^l = v_0^{l+1}$.
2. p_0 reads and p_1 reads: commutativity implies that $v_0^l = v_0^{l+1}$.
3. p_0 writes and p_1 reads: $v_0^l = v_0^{l+1}$, since the state of p_0 is the same in $D_l 0$ and D_{l+1} .

4. p_0 reads and p_1 writes: $v_1^l = v_1^{l+1}$, since the state of p_1 is the same in D_l and D_{l+1} .

By symmetric arguments we can show that for all l , $1 \leq l < k$, either $u_0^l = u_0^{l+1}$ or $u_1^l = u_1^{l+1}$. In a similar manner we show that either $v_1^l = u_1^l$ or $v_0^l = u_0^l$, by case analysis, depending on the type of operation taken in p_0 's step from C to C_l and in p_1 's step from C to C_l' :

1. p_0 writes and p_1 writes: commutativity implies that $v_0^l = v_0^{l+1}$.
2. p_0 reads and p_1 reads: commutativity implies that $v_0^l = v_0^{l+1}$.
3. p_0 writes and p_1 reads: $v_0^l = v_0^{l+1}$, since the state of p_0 is the same in D_l and D_l' .
4. p_0 reads and p_1 writes: $v_1^l = v_1^{l+1}$, since the state of p_1 is the same in D_l and D_l' .

Thus, either there exists some l such that $|v_0^l - v_1^l| \geq \frac{1}{2k} |pref_0(C) - pref_1(C)|$, or there exists some l such that, $|u_0^l - u_1^l| \geq \frac{1}{2k} |pref_0(C) - pref_1(C)|$. In the first case, the claim follows by taking $\sigma' = 0^l 1$, in the second case, the claim follows by taking $\sigma' = 1^l 0$. ■

Note that the validity condition implies that if p_i 's input in an initial configuration C is v_i then $pref_i(C) = v_i$. Starting with this fact and applying Lemma 8.1 iteratively, we can bound the rate at which a k -bounded algorithm converges. We get:

Theorem 8.2 *Let A be a k -bounded wait-free algorithm for approximate agreement between 2 processes. Then there exists an execution of A where processes start with inputs $\langle x_0, x_1 \rangle$ in which the time complexity is at least $\Omega(\log_{2k} \frac{|x_0 - x_1|}{\epsilon})$.*

Proof: Let C be an initial configuration in which processes have inputs $\langle x_0, x_1 \rangle$. We construct, inductively, a schedule σ_l such that σ_l is a sequence of l blocks and for $C_l = C\sigma_l$,

$$|pref_0(C_l) - pref_1(C_l)| \geq \left(\frac{1}{2k}\right)^l |pref_0(C) - pref_1(C)|.$$

This is done by applying Lemma 8.1. We have that $time(\sigma_l) = l$, since σ_l consists of l blocks. The validity condition implies that $pref_i(C) = x_i$. Thus, $|pref_0(C) - pref_1(C)| = |x_0 - x_1|$. The claim follows by noticing that it cannot be that both p_0 and p_1 have decided in a configuration D if $|pref_0(D) - pref_1(D)| > \epsilon$. ■

Remark 8.1 The case analysis in the proof of Lemma 8.1 can be extended to handle multi-writer multi-reader registers; thus, the above tradeoff applies also to algorithms that use multi-writer multi-reader atomic registers.

9 Properties of the Bias Function

In this section the interested reader may find the long postponed proofs of Lemma 4.1 through 4.4.

We begin with the rather straightforward proof of Lemma 4.1.

Lemma 4.1 *Let c^0, c^1 be nonnegative integers, and v^0, v^1, ε be real numbers, with $\varepsilon > 0$. Then $\text{bias}(v^0, v^1, c^0, c^1, \varepsilon) \in \text{range}(\{v^0, v^1\})$.*

Proof: Let $y = \text{bias}(v^0, v^1, c^0, c^1, \varepsilon)$. The claim is trivial if y is calculated in Line 1. Suppose y is calculated in Line 2. (The case where y is calculated in Line 3 is symmetric.) Then $y = v^1 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - \min\{c^1\varepsilon, |v^1|\})$. If the min is attained in the second term, then $y = v^1$ and the claim follows. So assume $c^1\varepsilon \leq v^1$, so $y = v^1 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - c^1\varepsilon)$. Assume $v^1 \geq v^0$. (A symmetric argument applies when $v^1 < v^0$.) Then $v^0 - v^1 \leq 0$, and clearly $y \leq v^1$. Since $|\frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - c^1\varepsilon)| \leq v^1 - v^0$, it follows that $y \geq v^0$. ■

The following is the proof of Lemma 4.2.

Lemma 4.2 *Let c^0, c^1 be nonnegative integers, and v^0, v^1, ε, m be real numbers, $\varepsilon > 0$, $m \geq 0$.*

(1) *Suppose $c^1 > c^0$ and $|v^1|/\varepsilon - m \leq c^1$. Then $|\text{bias}(v^0, v^1, c^0, c^1, \varepsilon) - v^1| \leq m\varepsilon$.*

(2) *Suppose $c^0 \geq c^1$ and $|v^0|/\varepsilon - m \leq c^0$. Then $|\text{bias}(v^0, v^1, c^0, c^1, \varepsilon) - v^0| \leq m\varepsilon$.*

Proof: We present the proof only for (2), the proof for (1) follows from symmetric arguments. Let $y = \text{bias}(v^0, v^1, c^0, c^1, \varepsilon)$. If y is calculated in Line 1 of bias , then $y = 0$ and $v^0 = 0$ and the claim follows. Hence, since $c^0 \geq c^1$ it follows that y is calculated in Line 3 of bias , i.e.,

$$y = v^0 + \frac{v^1 - v^0}{|v^0| + |v^1|}(|v^0| - \min\{c^0\varepsilon, |v^0|\}) .$$

If the min attains its value in the second term then $y = v^0$, and the claim follows. Otherwise, $c^0\varepsilon \leq |v^0|$; thus,

$$\begin{aligned} |y - v^0| &= |v^0 + \frac{v^1 - v^0}{|v^0| + |v^1|}(|v^0| - c^0\varepsilon) - v^0| \\ &= |\frac{v^1 - v^0}{|v^0| + |v^1|}(|v^0| - c^0\varepsilon)| \\ &= \frac{|v^1 - v^0|}{|v^0| + |v^1|}||v^0| - c^0\varepsilon| \\ &\leq ||v^0| - c^0\varepsilon| = |v^0| - c^0\varepsilon \leq m\varepsilon , \\ &\quad \text{by the hypothesis of the lemma.} \end{aligned}$$

■

Next is the proof of Lemma 4.3.

Lemma 4.3 *Let $c_0^0, c_0^1, c_1^0, c_1^1$ be nonnegative integers, and v^0, v^1, ε, m be real numbers, $\varepsilon > 0$ and $m \geq 0$. Suppose $\min\{c_0^0, c_0^1\} = \min\{c_1^0, c_1^1\} = 0$ and $|c_0^0 - c_1^0| + |c_0^1 - c_1^1| \leq m$. Then*

$$|\text{bias}(v^0, v^1, c_0^0, c_0^1, \varepsilon) - \text{bias}(v^0, v^1, c_1^0, c_1^1, \varepsilon)| \leq m\varepsilon.$$

Proof: Let $y_0 = \text{bias}(v^0, v^1, c_0^0, c_0^1, \varepsilon)$, and $y_1 = \text{bias}(v^0, v^1, c_1^0, c_1^1, \varepsilon)$.

If $v^0 = v^1 = 0$ then both y_0 and y_1 are calculated in Line 1 of *bias*, i.e., $y_0 = y_1 = 0$ and the claim follows.

Now assume y_0 is calculated in Line 2 of *bias*, while y_1 is calculated in Line 3 of *bias* (the reverse case is symmetric). Thus, $c_0^0 < c_0^1$, while $c_1^1 \leq c_1^0$. Thus, by assumption, $c_0^0 = c_1^1 = 0$. Since $|c_0^0 - c_1^0| + |c_0^1 - c_1^1| \leq m$, it follows that $|c_1^0| + |c_0^1| = c_1^0 + c_0^1 \leq m$. Thus, $\min\{c_1^0, |v^0|/\varepsilon\} + \min\{c_0^1, |v^1|/\varepsilon\} \leq m$. So, $\min\{c_1^0\varepsilon, |v^0|\} + \min\{c_0^1\varepsilon, |v^1|\} \leq m\varepsilon$. We have

$$y_0 = v^1 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - \min\{c_0^1\varepsilon, |v^1|\}) \quad \text{and} \quad y_1 = v^0 + \frac{v^1 - v^0}{|v^0| + |v^1|}(|v^0| - \min\{c_1^0\varepsilon, |v^0|\}).$$

Thus,

$$\begin{aligned} |y_0 - y_1| &= |v^1 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - \min\{c_0^1\varepsilon, |v^1|\}) - v^0 - \frac{v^1 - v^0}{|v^0| + |v^1|}(|v^0| - \min\{c_1^0\varepsilon, |v^0|\})| \\ &= |v^1 - v^0 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^0| + |v^1|) - \frac{v^0 - v^1}{|v^0| + |v^1|}(\min\{c_0^1\varepsilon, |v^1|\} + \min\{c_1^0\varepsilon, |v^0|\})| \\ &= \frac{|v^0 - v^1|}{|v^0| + |v^1|} |\min\{c_0^1\varepsilon, |v^1|\} + \min\{c_1^0\varepsilon, |v^0|\}| \\ &\leq |\min\{c_0^1\varepsilon, |v^1|\} + \min\{c_1^0\varepsilon, |v^0|\}| = \min\{c_0^1\varepsilon, |v^1|\} + \min\{c_1^0\varepsilon, |v^0|\} \leq m\varepsilon, \end{aligned}$$

as needed.

Now assume that both y_0 and y_1 are calculated in Line 2 of *bias* (the case where y_0 and y_1 are calculated in Line 3 of *bias* is symmetric), i.e.,

$$y_0 = v^1 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - \min\{c_0^1\varepsilon, |v^1|\}) \quad \text{and} \quad y_1 = v^1 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - \min\{c_1^1\varepsilon, |v^1|\}).$$

If for y_0 the min is attained in the second term, then $c_0^1\varepsilon \geq |v^1|$, and $y_0 = v^1$; since $|c_0^1 - c_1^1| \leq m$ it follows that $c_1^1 \geq |v^1|/\varepsilon - m$. Because y_1 is calculated in Line 2, $c_1^0 < c_1^1$ and the claim follows from Lemma 4.2 (1). A similar argument applies if for y_1 the min is attained in the second term. So assume that for both y_0 and y_1 the min is attained in the first term. Thus,

$$|y_0 - y_1| = |v^1 + \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - c_0^1\varepsilon) - v^1 - \frac{v^0 - v^1}{|v^0| + |v^1|}(|v^1| - c_1^1\varepsilon)|$$

$$\begin{aligned}
&= \left| \frac{v^0 - v^1}{|v^0| + |v^1|} (c_1^1 \varepsilon - c_0^1 \varepsilon) \right| \\
&= \frac{|v^0 - v^1|}{|v^0| + |v^1|} |(c_1^1 \varepsilon - c_0^1 \varepsilon)| \\
&\leq |(c_1^1 \varepsilon - c_0^1 \varepsilon)| = \varepsilon |c_1^1 - c_0^1| \leq m \varepsilon,
\end{aligned}$$

as needed. ■

In the proof of the next lemma we use the following two facts:

Claim 9.1 *If x, y, x', y' are real numbers, and for some δ , $|x - x'| \leq \delta$ and $|y - y'| \leq \delta$, then $\left| \frac{|x|(y-x)}{|x|+|y|} - \frac{|x'|(y'-x')}{|x'|+|y'|} \right| \leq 3\delta$.*

We prove this claim by first showing that $\left| \frac{x(y-x)}{x+y} - \frac{x'(y'-x')}{x'+y'} \right| \leq 3\delta$, using calculus, then handling the absolute values by case analysis.

Claim 9.2 *If x, y, x', y' are real numbers, and for some δ , $|x - x'| \leq \delta$ and $|y - y'| \leq \delta$, then $\left| \frac{(y-x)}{|x|+|y|} - \frac{(y'-x')}{|x'|+|y'|} \right| \leq \frac{2\delta}{\min(|x|+|y|, |x'|+|y'|)}$.*

We prove this claim by straightforward calculations and a case analysis. Finally, we can now prove Lemma 4.4.

Lemma 4.4 *Let c^0, c^1 be nonnegative integers, and $v_0^0, v_0^1, v_1^0, v_1^1, \varepsilon, \delta$ be real numbers, with $\varepsilon > 0$, $\delta \geq 0$. Suppose $|v_0^0 - v_1^0| \leq \delta$ and $|v_0^1 - v_1^1| \leq \delta$. Then*

$$|\text{bias}(v_0^0, v_0^1, c^0, c^1, \varepsilon) - \text{bias}(v_1^0, v_1^1, c^0, c^1, \varepsilon)| \leq 6\delta.$$

Proof: Let $y_0 = \text{bias}(v_0^0, v_0^1, c^0, c^1, \varepsilon)$, and $y_1 = \text{bias}(v_1^0, v_1^1, c^0, c^1, \varepsilon)$. If $v_0^0 = v_0^1 = 0$ then $y_0 = 0$. Thus, $|v_1^0| \leq \delta$ and $|v_1^1| \leq \delta$. So from Lemma 4.1 it follows that $|y_1| \leq \delta$ and the claim follows. The case $v_1^0 = v_1^1 = 0$ follows from symmetric arguments. So assume at least one of v_0^0, v_0^1 is nonzero and similarly for at least one of v_1^0, v_1^1 .

Assume that $c^0 < c^1$, i.e., y_0 and y_1 are calculated in Line 2. (The other case, where $c^1 \leq c^0$ and y_0 and y_1 are calculated in Line 3, is symmetric.) Then

$$y_0 = v_0^1 + \frac{v_0^0 - v_0^1}{|v_0^0| + |v_0^1|} (|v_0^1| - \min\{c^1 \varepsilon, |v_0^1|\}) \quad \text{and} \quad y_1 = v_1^1 + \frac{v_1^0 - v_1^1}{|v_1^0| + |v_1^1|} (|v_1^1| - \min\{c^1 \varepsilon, |v_1^1|\}).$$

First, assume the min for y_0 is attained in the second term; then $y_0 = v_0^1$. In this case, if the min for y_1 is also attained in the second term, then $y_1 = v_1^1$, and the claim follows. On the other hand, suppose the min for y_1 is attained in the first term. Since the min for y_0 is

attained in the second term, $c^1 \varepsilon \geq |v_0^1| \geq |v_1^1| - \delta$. Applying Lemma 4.2 (1) with $m = \delta/\varepsilon$, we get that $|y_1 - v_1^1| \leq \delta$. Since $|v_0^1 - v_1^1| \leq \delta$, we have $|y_0 - y_1| \leq 2\delta$.

Now assume that in both cases the min is attained in the first term. In particular, $c^1 \varepsilon \leq |v_1^1|$ and $c^1 \varepsilon \leq |v_0^1|$. We have,

$$\begin{aligned}
|y_0 - y_1| &= |v_0^1 + \frac{v_0^0 - v_0^1}{|v_0^0| + |v_0^1|}(|v_0^1| - c^1 \varepsilon) - v_1^1 - \frac{v_1^0 - v_1^1}{|v_1^0| + |v_1^1|}(|v_1^1| - c^1 \varepsilon)| \\
&\leq |v_0^1 - v_1^1| + |\frac{v_0^0 - v_0^1}{|v_0^0| + |v_0^1|}(|v_0^1| - c^1 \varepsilon) - \frac{v_1^0 - v_1^1}{|v_1^0| + |v_1^1|}(|v_1^1| - c^1 \varepsilon)| \\
&\leq \delta + |\frac{v_0^0 - v_0^1}{|v_0^0| + |v_0^1|}(|v_0^1| - c^1 \varepsilon) - \frac{v_1^0 - v_1^1}{|v_1^0| + |v_1^1|}(|v_1^1| - c^1 \varepsilon)| \\
&\leq \delta + |\frac{|v_0^1|(v_0^0 - v_0^1)}{|v_0^0| + |v_0^1|} - \frac{|v_1^1|(v_1^0 - v_1^1)}{|v_1^0| + |v_1^1|}| + |\frac{v_0^0 - v_0^1}{|v_0^0| + |v_0^1|}c^1 \varepsilon - \frac{v_1^0 - v_1^1}{|v_1^0| + |v_1^1|}c^1 \varepsilon| \\
&\leq 4\delta + c^1 \varepsilon |\frac{v_0^0 - v_0^1}{|v_0^0| + |v_0^1|} - \frac{v_1^0 - v_1^1}{|v_1^0| + |v_1^1|}|, \quad \text{by Claim 9.1,} \\
&\leq 4\delta + c^1 \varepsilon \frac{2\delta}{\min(|v_0^1| + |v_0^0|, |v_1^1| + |v_1^0|)}, \quad \text{by Claim 9.2,} \\
&\leq 4\delta + c^1 \varepsilon \frac{2\delta}{\min(|v_0^1|, |v_1^1|)} \leq 4\delta + c^1 \varepsilon \frac{2\delta}{c^1 \varepsilon} \leq 6\delta.
\end{aligned}$$

■

10 Discussion and Further Research

For approximate agreement, the answer to the question whether wait-free algorithms are fast is not binary, rather it is quantitative: we have presented a relatively fast, $O(\log n)$ time, wait-free algorithm for n -process approximate agreement. On the other hand, $\log n$ is a lower bound on the time complexity of any wait-free approximate agreement algorithm, and there exists an $O(1)$ time non-wait-free algorithm.

Using the emulators of [5], our algorithms can be translated into algorithms that work in message-passing systems. The algorithms have the same time complexity (in complete networks) and are resilient to the failure of a majority of the processes.

There are many ways in which our work can be extended. An interesting direction is to consider the impact on our results of using other shared memory primitives. For example, if powerful *Read-Modify-Write* registers are used, then a constant time wait-free approximate agreement algorithm can be devised. What happens if *multi-writer* multi-reader registers are used? The existence of faster wait-free algorithms using these primitives will imply a lower bound on the *time complexity* (in normal executions) of any implementation of multi-writer registers from single-writer registers.

Another avenue of research is to see whether the techniques presented in this paper, both for algorithms and lower bounds, can be applied to other problems. We believe, for example, that the $O(1)$ time algorithm for 2-process approximate agreement can be generalized to *any* decision problem of size 2, using the characterization result of [8]. It is interesting to explore whether similar results can be proved for problems that require repeated coordination (e.g., ℓ -exclusion).

Finally, there remains the fundamental unanswered question raised by this work: Can wait-free (highly resilient) computation be performed at the price of no more than a logarithmic slowdown? Even more strongly, are there $O(\log n)$ time wait-free algorithms for *all* problems that have wait-free solutions?

Following a preliminary version of our work, first steps were made towards answering this question in the context of randomized computation [46]. Based on the alternated-interleaving method presented in Section 6.2, it is shown that any *decision problem* that has a wait-free or expected wait-free¹⁰ solution algorithm, has an expected wait-free algorithm with the same worst case time complexity, that takes only $O(\log n)$ expected time¹¹ in fault-free executions. However, the above question itself is still far from being answered.

Acknowledgements:

We would like to thank Jennifer Welch for careful reading of an earlier version of the paper and for many helpful comments. Thanks are also due to Cynthia Dwork, Maurice Herlihy, Mike Saks, Marc Snir and Heather Woll, for helpful discussions on the topic of this paper.

¹⁰An expected wait-free algorithm is a randomized algorithm that is only expected, rather than guaranteed, to terminated within a finite number of steps.

¹¹This is optimal by a straightforward extension of our lower bound to the case of randomized computation (see [46]).

References

- [1] Y. Afek, H. Attiya, D. Dolev, E. Gafni, M. Merritt and N. Shavit, "Atomic Snapshots of Shared Memory," *Proc. 9th ACM Symp. on Principles of Distributed Computing, Quebec-City*, August 1990, pp. 1-14.
- [2] J. Anderson, "Composite Registers," *Proc. 9th ACM Symp. on Principles of Distributed Computing, Quebec-City*, August 1990 pp. 15-30.
- [3] E. Arjomandi, M. Fischer and N. Lynch, "Efficiency of Synchronous Versus Asynchronous Distributed Systems," *Journal of the ACM*, Vol. 30, No. 3 (1983), pp. 449-456.
- [4] J. Aspnes and M. Herlihy, "Fast Randomized consensus Using Shared Memory," *Journal of Algorithms*, Vol. 11, pp. 441-461, September 1990.
- [5] H. Attiya, A. Bar-Noy and D. Dolev, "Sharing Memory Robustly in Message-Passing Systems," *9th Annual ACM Symposium on Principles of Distributed Computing (PODC)*, Quebec-City, August 1990, pp. 363-376.
Expanded version: Technical Memo MIT/LCS/TM-423, Laboratory for Computer Science, MIT, February 1990.
- [6] H. Attiya and N. Lynch, "Time Bounds for Real-Time Process Control in the Presence of Timing Uncertainty," in *proceedings of the 10th IEEE Real-Time Systems Symposium*, Santa-Monica, December 1989, pp. 268-284.
Expanded version: Technical Memo MIT/LCS/TM-403, Laboratory for Computer Science, MIT, July 1989.
- [7] H. Attiya, N. Lynch and N. Shavit, "Are Wait-Free Algorithms Fast?" *31st Annual Symposium on the Foundations of Computer Science, St. Louis*, October 1990.
- [8] O. Biran, S. Moran and S. Zaks, "A Combinatorial Characterization of the Distributed Tasks which are Solvable in the Presence of One Faulty Processor," *Journal of Algorithms*, Vol. 11, pp. 420-440, September 1990.
- [9] B. Coan and C. Dwork, "Simultaneity is Harder than Agreement", *Proc. 5th IEEE Symposium on Reliability in Distributed Software and Database Systems*, pp. 141-150, 1986.
- [10] C. Dwork and D. Skeen, "The Inherent Cost of Nonblocking Commitment," *Proc. 2nd ACM Symp. on Principles of Distributed Computing*, 1983, pp. 1-11.
- [11] S. Cook, C. Dwork and R. Reischuk, "Upper and Lower Time Bounds for Parallel RAMS Without Simultaneous Writes," *SIAM J. Computing*, Vol. 15, No. 1, 1986, pp. 87-98.
- [12] R. Cole and O. Zajicek, "The APRAM: Incorporating Asynchrony into the PRAM model," *Proc. 1st ACM Symp. on Parallel Algorithms and Architectures*, 1989, pp. 169-178.

- [13] R. Cole and O. Zajicek, "The Expected Advantage of Asynchrony," *Proc. 2nd ACM Symp. on Parallel Algorithms and Architectures*, 1990, pp. 85-94.
- [14] D. Dolev, E. Gafni and N. Shavit, "Toward a Non-Atomic Era: ℓ -Exclusion as a Test Case," *Proc. 20th ACM Symp. on the Theory of Computing*, 1988, pp. 78-92.
- [15] D. Dolev, N. Lynch, S. Pinter, E. Stark and W. Weihl, "Reaching Approximate Agreement in the Presence of Faults," *Journal of the ACM*, Vol. 33, No. 3, 1986, pp. 499-516.
- [16] D. Dolev, R. Reischuk and H. R. Strong, "Eventual Is Earlier Than Immediate," *Proc. 23rd IEEE Symp. on Foundations of Computer Science*, 1982, pp. 196-203.
- [17] D. Dolev, C. Dwork and L. Stockmeyer, "On the Minimal Synchrony Needed for Distributed Consensus," *Journal of the ACM*, Vol. 34, No. 1 (January 1987), pp. 77-97.
- [18] C. Dwork and Y. Moses, "Knowledge and Common Knowledge in a Byzantine Environment: Crash Failures," to appear in *Information and Computation*.
- [19] A. Fekete, "Asymptotically Optimal Algorithms for Approximate Agreement," *Proc. 5th ACM Symp. on Principles of Distributed Computing*, 1986, pp. 73-87.
- [20] A. Fekete, "Asynchronous Approximate Agreement," *Proc. 6th ACM Symp. on Principles of Distributed Computing*, 1987, pp. 64-76.
- [21] M. Fischer, N. Lynch and M. Paterson, "Impossibility of Distributed Consensus with One Faulty Processor," *Journal of the ACM*, Vol. 32, No. 2 (1985), pp. 374-382.
- [22] P. Gibbons, "Towards Better Shared Memory Programming Models," *Proc. 1st ACM Symp. on Parallel Algorithms and Architectures*, 1989, pp. 169-178.
- [23] M. P. Herlihy, "Impossibility and Universality Results for Wait-Free Synchronization," *Proc. 7th ACM Symp. on Principles of Distributed Computing*, 1988, pp. 276-290.
- [24] P. Kanellakis and A. Shvartsman, "Efficient Parallel Algorithms can be Made Robust," *Proc. 8th ACM Symp. on Principles of Distributed Computing*, 1989, pp. 211-221.
- [25] Z. Kedem, K. Palem and P. Spirakis, "Efficient Robust Parallel Computations," *Proc. 22nd ACM Symp. on Theory of Computing*, 1990, pp. 138-148.
- [26] L. Lamport, "The Synchronization of Independent Processes," *Acta Informatica*, Vol. 7, No. 1 (1976), pp. 15-34.
- [27] L. Lamport, "Proving the Correctness of Multiprocess Programs," *IEEE Transactions on Software Engineering*, Vol. SE-3, No. 2 (March 1977) pp. 125-143.
- [28] L. Lamport, "On Interprocess Communication. Part I: Basic Formalism," *Distributed Computing* 1, 2 1986, 77-85.

- [29] L. Lamport, "On Interprocess Communication. Part II: Algorithms," *Distributed Computing* 1, 2 1986, pp. 86-101.
- [30] L. Lamport, R. Shostak and M. Pease, "The Byzantine Generals Problem," *ACM Transactions on Programming Languages and Systems*, Vol. 4, No. 3 (July 1982), pp. 382-401.
- [31] B. Lampson, "Hints for Computer System Design", in *Proc. 9th ACM Symposium on Operating Systems Principles*, 1983, pp. 33-48.
- [32] M. Li, J. Tromp and P. M.B. Vitanyi, "How to Share Concurrent Wait-Free Variables," *ICALP 1989*. Expanded version: Report CS-R8916, CWI, Amsterdam, April 1989.
- [33] M. Loui and H. Abu-Amara, "Memory Requirements for Agreement Among Unreliable Asynchronous Processes," *Advances in Computing Research*, Vol. 4, JAI Press, Inc., 1987, 163-183.
- [34] N. Lynch and M. Fischer, "On Describing the Behavior and Implementation of Distributed Systems," *Theoretical Computer Science*, Vol. 13, No. 1 (January 1981), pp. 17-43.
- [35] N. Lynch and K. Goldman, *Lecture notes for 6.852*. MIT/LCS/RSS-5, Laboratory for Computer Science, MIT, 1989.
- [36] S. Mahaney and F. Schneider, "Inexact Agreement: Accuracy, Precision, and Graceful Degradation," *Proc. 4th ACM Symp. on Principles of Distributed Computing*, 1985, pp. 237-249.
- [37] C. Martel, A. Park and R. Subramonian, "Optimal Asynchronous Algorithms for Shared Memory Parallel Computers," Technical Report CSE-89-8, Division of Computer Science, University of California, Davis, July 1989.
- [38] C. Martel, R. Subramonian and A. Park, "Asynchronous PRAMs are (Almost) as Good as Synchronous PRAMs," *Proc. 31st IEEE Symp. on Foundations of Computer Science*, 1990, pp. 590-599.
- [39] M. Merritt, F. Modugno and M. Tuttle, "Time Constrained Automata," manuscript, November 1988.
- [40] Y. Moses and M. Tuttle, "Programming Simultaneous Actions using Common Knowledge," *Algoritmica*, Vol. 3, 1988, pp. 121-169.
- [41] N. Nishimura, "Asynchronous Shared Memory Parallel Computation," *Proc. 2nd ACM Symp. on Parallel Algorithms and Architectures*, pp. 76-84, 1990.
- [42] G. Peterson, "Concurrent Reading While Writing," *ACM Transactions on Programming Languages and Systems*, Vol. 5, No. 1 (January 1983), pp. 46-55.
- [43] G. Peterson, and J. Burns, "Concurrent Reading While Writing II : The Multi-Writer Case," *Proc. 28th IEEE Symp. on Foundations of Computer Science*, 1987, pp. 383-392.

- [44] G. Peterson and M. Fischer, "Economical Solutions for the Critical Section Problem in a Distributed System," *Proc. 9th ACM Symp. on Theory of Computing*, 1977, pp. 91-97.
- [45] R. Schaffer, "On the Correctness of Atomic Multi-Writer Registers," MIT/LCS/TM-364, June 1988.
- [46] M. Saks, N. Shavit and H. Woll, "Optimal Time Randomized Consensus - Making Resilient Algorithms Fast in Practice," *Proc. of the 2nd ACM Symposium on Discrete Algorithms*, pp. 351-362 January 1991.
- [47] D. Skeen, "Crash Recovery in a Distributed Database System," Memorandum No. UCB/ERL M82/45, Electronics Research Laboratory, Berkeley, May 1982.
- [48] P. Vitanyi and B. Awerbuch, "Atomic Shared Register Access by Asynchronous Hardware," *Proc. 27th IEEE Symp. on Foundations of Computer Science*, pp. 233-243, 1986.
- [49] J. Wyllie, *The Complexity of Parallel Computation*, Ph.D. thesis, Cornell University, August 1979. Technical Report TR 79-387, Department of Computer Science.

OFFICIAL DISTRIBUTION LIST

DIRECTOR Information Processing Techniques Office Defense Advanced Research Projects Agency (DARPA) 1400 Wilson Boulevard Arlington, VA 22209	2 copies
OFFICE OF NAVAL RESEARCH 800 North Quincy Street Arlington, VA 22217 Attn: Dr. Gary Koop, Code 433	2 copies
DIRECTOR, CODE 2627 Naval Research Laboratory Washington, DC 20375	6 copies
DEFENSE TECHNICAL INFORMATION CENTER Cameron Station Alexandria, VA 22314	12 copies
NATIONAL SCIENCE FOUNDATION Office of Computing Activities 1800 G. Street, N.W. Washington, DC 20550 Attn: Program Director	2 copies
HEAD, CODE 38 Research Department Naval Weapons Center China Lake, CA 93555	1 copy